



Norwegian University
of Life Sciences

Master's Thesis 2023 30 ECTS
Faculty of Science and Technology

Satellite based methane emission estimation for flaring activities in oil and gas industry: A data-driven approach(SMEEF-OGI)

Muhammad Uzair Aftab
Data Science

PREFACE

This thesis signifies the culmination of my master's journey in data science at the Norwegian University of Life Sciences. Over the course of these transformative years at NMBU, I have been privileged to cross paths with a host of exceptional individuals, to whom I wish to extend my sincere gratitude.

Primary among these is Eik Lab, for providing an platform for learning beyond textbooks, helping me delve deeper into my interests. Without you, I'd just be just your average Data Scientist.

I am immensely grateful to Kristian Omberg, whose guidance and wisdom were instrumental throughout the course of this thesis. His patience in navigating my thesis was truly impressive, and I'm thankful for it.

Equally, Ola Omberg deserves special mention for his invaluable advice, and the perseverance he showed in waiting three long grueling years before offering his first commendation. As a wise man once said, "There's no inflation on atta boys."

My heartfelt thanks go to Madelen and Mari for welcoming me into their circle and allowing their duo to evolve into a trio. Similarly, I owe a debt of gratitude to Marthe-Andrea, who played an irreplaceable role in altering my worldview and softening my judgments. And for making my life a little less boring, and a bit more social.

Aleksander Eriksen, my partner in our entrepreneurial venture, Njord, deserves special recognition for his unwavering faith in me. The late-night coding sessions fueled by caffeine, takeaways, and passionate rants have been as educational as they were enjoyable. And for the mishap with the Apple watch, my apologies once again.

Thank you to Sander Wittwer, and Balder Klanderud, for serving as my de-facto Wikipedia for all things software, Linux, and Rust. In Rust we trust.

A special mention goes to Emil Skaar, who trusted me with my first industry project, instilling in me the discipline to scrutinize my code at least 22 times before committing it to Git.

To Lavanyan, thank you for providing the opportunity to enhance my leadership skills as your mentor. I hope my push for perfection wasn't too demanding. Likewise, a big thank you to Trym, my 'work wife' for the past two years. I appreciate your tolerance of my procrastination throughout our late nights of work.

To both Lavanyan, and Trym, I'm looking forward to see what you guys are going to achieve in your life.

Last, but not least.

Thank you, Kristian, and Ola, for having some faith in an eager, over-confident 18-year-old, who experienced the Dunning-Kruger effect early in his life.



Muhammad Uzair Aftab
15th of May 2023, Ås

EXECUTIVE SUMMARY

Climate change, precipitated in part by greenhouse gas emissions, presents a critical global challenge. Methane, a highly potent greenhouse gas with a global warming potential of 80 times that of carbon dioxide, is a significant contributor to this crisis. Sources of methane emissions include the oil and gas industry, agriculture, and waste management, with flaring in the oil and gas industry constituting a significant emission source.

Flaring, a standard process in the Oil and gas industry is often assumed to be 98% efficient when converting methane to less harmful carbon dioxide. However, recent research from the University of Michigan, Stanford, the Environmental Defense Fund, and Scientific Aviation indicates that the widely accepted 98% efficiency of flaring in converting methane to carbon dioxide, a less harmful greenhouse gas, may be inaccurate. This investigation reevaluates the flaring process's efficiency and its role in methane conversion.

This work focuses on creating a method to independently calculate methane emissions from oil and gas activities to solve this issue. Satellite data, which is a helpful tool for calculating greenhouse gas emissions from various sources, is included in the suggested methodology. In addition to standard monitoring techniques, satellite data offers an independent, non-intrusive, affordable, and continuous monitoring approach.

Based on this, the problem statement for this work is the following

"How can a data-driven approach be developed to enhance the accuracy and quality of methane emission estimation from flaring activities in the Oil and Gas industry, using satellite data from selected platforms to detect and quantify future emissions based on Machine learning more effectively?"

To achieve this, the following objectives and activities were performed.

- Theoretical Framework and key concepts
- Technical review of the current state-of-the-art satellite platforms and existing literature.
- Development of a Proof of Concept
- Proposing an evaluation of the method
- Recommendations and further work

This work has adopted a systematic approach, starting with a comprehensive theoretical framework to understand the utilization of flaring, the environmental implications of methane, the current state-of-the-art of research, and the state-of-the-art in the field of remote sensing via satellites.

Based upon the framework developed during the initial phases of this work, a data-driven methodology was formulated, utilizing the VIIRS dataset to get geographical areas of interest. Hyperspectral and methane data were aggregated from the Sentinel-2 and Sentinel-5P satellite dataset. This information was processed via a proposed pipeline, with initial alignment and enhancement. In this work, the images were enhanced by calculating the Normalized Burn Index.

The result was a dataset containing the location of known flare sites, with data from both the Sentinel-2, and the Sentinel-5P satellite.

The results underscore the disparities in coverage between Sentinel-2 and Sentinel-5P data, a factor that could potentially influence the precision of methane emission estimates. The applied preprocessing techniques markedly enhanced data clarity and usability, but their efficacy may hinge on the flaring sites' specific characteristics and the raw data quality. Moreover, despite certain limitations, the combination of Sentinel-2 and Sentinel-5P data effectively yielded a comprehensive dataset suitable for further analysis.

In conclusion, this project introduces an encouraging methodology for estimating methane emissions from flaring activities within the oil and gas industry. It lays a foundational steppingstone for future research, continually enhancing the precision and quality of data in combating climate change. This methodology can be seen in the flow chart below.

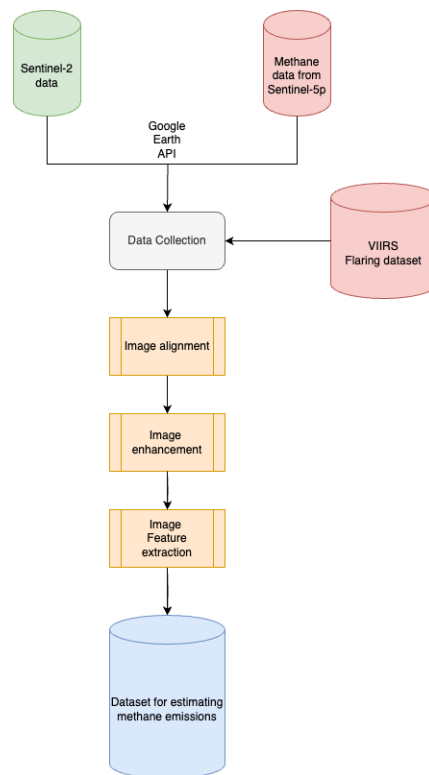


Figure 1 An overview of the proposed dataset generation method

Based on the work done in this project, future work could focus on incorporating alternative sources of methane data, broadening the areas of interest through industry collaboration, and attempting to extract further features through image segmentation methods. This project signifies a start, paving the way for subsequent explorations to build upon.

SAMMENDRAG

Klimaendringer, delvis utløst av klimagassutslipp, utgjør en kritisk global utfordring. Metan, en svært potent drivhusgass med et globalt oppvarmingspotensial på 80 ganger karbondioksid, er en betydelig bidragsyter til denne krisen. Kilder til metanutslipp inkluderer olje- og gassindustrien, landbruket og avfallshåndteringen, med faking i olje- og gassindustrien som en betydelig utslippskilde.

Faking, en standardprosess i olje- og gassindustrien, antas ofte å være 98 % effektiv ved omdannelse av metan til mindre skadelig karbondioksid. Nyere forskning fra University of Michigan, Stanford, Environmental Defense Fund og Scientific Aviation indikerer imidlertid at den allment aksepterte effektiviteten på 98 % av faking ved konvertering av metan til karbondioksid, en mindre skadelig klimagass, kan være unøyaktig. Denne undersøkelsen revurderer fakkelprosessens effektivitet og dens rolle i metankonvertering.

Dette arbeidet fokuserer på å lage en metode for uavhengig å beregne metanutslipp fra olje- og gassvirksomhet for å løse dette problemet. Satellittdata, som er et nyttig verktøy for å beregne klimagassutslipp fra ulike kilder, er inkludert i den foreslåtte metodikken. I tillegg til standard overvåkingsteknikker, tilbyr satellittdata en uavhengig, ikke-påtrengende, rimelig og kontinuerlig overvåkingstilnærming.

På bakgrunn av dette er problemstillingen for dette arbeidet følgende

"Hvordan kan en datadrevet tilnærming utvikles for å forbedre nøyaktigheten og kvaliteten på estimering av metanutslipp fra faklingsaktiviteter i olje- og gassindustrien, ved å bruke satellittdata fra utvalgte plattformer for å oppdage og kvantifisere fremtidige utslipp basert på maskinlæring mer effektivt?"

For å oppnå dette ble følgende mål og aktiviteter utført.

- Teoretisk rammeverk og sentrale begreper
- Teknisk gjennomgang av dagens toppmoderne satellittplattformer og eksisterende litteratur.
- Utvikling av et Proof of Concept
- Foreslå en evaluering av metoden
- Anbefalinger og videre arbeid

Dette arbeidet har tatt i bruk en systematisk tilnærming, som starter med et omfattende teoretisk rammeverk for å forstå bruken av faking, de miljømessige implikasjonene av metan, den nåværende «state-of-the-art» av forskning, og «state-of-the-art» i felt for fjernmåling via satellitter.

Basert på rammeverket utviklet i de innledende fasene av dette arbeidet, ble det formulert en datadrevet metodikk, som benytter VIIRS-datasettet for å få geografiske områder av interesse. Hyperspektrale data og metandata ble samlet fra Sentinel-2 og Sentinel-5P satellittdatasettet. Denne informasjonen ble behandlet via en foreslått rørledning, med innledende justering og forbedring. I dette arbeidet ble bildene forbedret ved å beregne den normaliserte brennindeksen.

Resultatet var et datasett som inneholdt plasseringen av kjente fakkellsteder, med data fra både Sentinel-2 og Sentinel-5P-satellitten.

Resultatene understreker forskjellene i dekningen mellom Sentinel-2- og Sentinel-5P-data, en faktor som potensielt kan påvirke nøyaktigheten av metanutslippsestimater. De anvendte forbehandlingsteknikkene forbedret dataklarheten og brukervennligheten markant, men deres effektivitet kan avhenge av fakkellstedenes spesifikke egenskaper og rådatakvaliteten. Dessuten, til tross for visse begrensninger, ga kombinasjonen av Sentinel-2 og Sentinel-5P-data effektivt et omfattende datasett egnet for videre analyse.

Avslutningsvis introduserer dette prosjektet en oppmuntrende metodikk for å estimere metanutslipp fra faking i olje- og gassindustrien. Den legger et grunnleggende springbrett for fremtidig forskning, og forbedrer kontinuerlig presisjonen og kvaliteten på data for å bekjempe klimaendringer. Denne metodikken kan sees i flytskjemaet nedenfor.

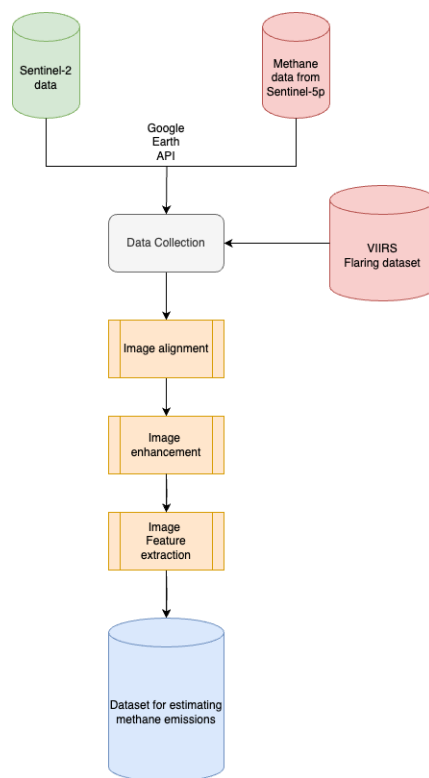


Figure 2 En oversikt over den foreslåtte datasettgenereringsmetoden

Basert på arbeidet som er gjort i dette prosjektet, kan fremtidig arbeid fokusere på å innlemme alternative kilder til metan data, utvide interesseområdene gjennom industrisamarbeid og forsøke å trekke ut ytterligere detaljer gjennom bilde-segenteringsmetoder. Dette prosjektet legger et grunnlag, og baner vei for påfølgende utforskninger å bygge videre på.

CONTENTS

List of figures.....	viii
List of Tables	ix
Abbreviations.....	xii
Introduction	1
Problem description.....	2
Limitations	3
Theory and key concepts	5
Methane and its environmental impact	5
What is flaring and the emission from flaring operations?	6
What are the market drivers for monitoring methane EMISSIONS?	7
Overview of satellite technology	9
The application of satellites in Methane emission monitoring.....	10
Image enhancement methods for Hyperspectral data.....	12
Spatial binning of satellite data	13
Introduction to Machine Learning	14
What is Machine Learning?	14
Methods of supervised machine learning	18
Data Quality and model assessment	22
Satellite Data Quality concepts and requirements	24
Satellite Imagery: Resolution vs. Accuracy	24
Spectral range	25
Signal-TO-Noise ratio (SNR)	25
Application User Interface	26
Google Earth Engine API	26
Copernicus Open Access HUB.....	26
Summary	27
Technical review	29
Flaring monitoring – how is emission from Flaring currently being MONITORED?	29
Emission monitoring via satellites.....	30
Technical Review of available satellite platforms	31
Sentinel-5 and 2	31
Landsat-8	33
WorldView-3	35
GHGsat- MethaneSAT	37
Visible Infrared Imaging Radiometer Suite (VIIRS).....	38
Literature review of relevant research	40

General improvements In data quality	40
Improvements based on Machine Learning and existing datasets.....	44
Specific work on THE development of data-driven methods for improving methane emission from satellites	45
Summary	49
Method - Development of data -preprocessing	51
Review of available data sources	51
Sentinel-2	51
Sentinel-5 - Tropomi (TROPOspheric Monitoring Instrument).....	54
Google Earth Engine	57
VIIRS dataset	58
Data preprocessing method – Description	60
Development of an «Proof of Concept”	62
Data acquisition and preprocessing.....	62
Result – presentation of the POC.....	65
Evaluation of the method	69
Testplan for future testing	69
Objectives	69
Test plan components	69
Discussion	71
Coverage differences between Sentinel-2 and Sentinel-5p.....	71
Preprocessing TECHNIQUES and their effectiveness	71
Application of the NBR Index.....	71
Integration of Sentinel-2 and Sentinel-5P data	71
Recommendations and further work.....	72
Further work	72
Image Segmentation.....	72
Feature Extraction	72
Data Augmentation.....	72
Exploration of Other Satellite Data Sources	73
Integration of other Ground Truth Data	73
Bibliography	74
Attachments	80
Sentinel 5P methane specifications	80
Flaring Monitor – User Interface	81

List of figures	iii
Figure 1 An overview of the proposed dataset generation method	iii
Figure 2 En oversikt over den foreslåtte datasettgenereringsmetoden	v
Figure 3 Major sources and sinks for methane. The size of the arrow indicates the relative contribution a source makes to the global total. Methane's lifetime is about nine years before oxidizing agents convert it into carbon dioxide. Image from [2]	1
Figure 4 An overview of methane emissions from fossil fuel exploitation [7]	5
Figure 5 Overview of typical gas Flare system	6
Figure 6 Example of a satellite in orbit around the Earth. Taken from [16]	9
Figure 7 Overview of a typical satellite Payload. Image taken from [18].	10
Figure 8 Absorption spectrum of methane (green), water (black), carbon dioxide (red), and ethane (blue) in the infrared region of the electromagnetic spectrum [19].	11
Figure 9 Three primary types of machine learning. Image taken from [20].	14
Figure 10 A cost function - In the above image the farther the points is from the straight red line higher the error in predicted value w.r.t ground truth value. Figure taken from [70].	15
Figure 11 General representation of a reinforcement learning scenario [20].	17
Figure 12 Conceptual overview of a neural network [23].	19
Figure 13 A conceptual overview of a single neuron [23].	20
Figure 14 An architectural overview of a convolutional neural network. Taken from [69].	21
Figure 15 Conceptual design of a methane emission tool	27
Figure 16 The operating principle of a transit-time-based ultrasonic flowmeter [30].	29
Figure 17 Example of a satellite application is CO2 measurement from orbit, where the satellite is scanning the earth's surface via its sensory payload [31].	30
Figure 18 Example of Sentinel-5P NO2 data visualized [33]	31
Figure 19 Acquired on 27 June 2015 at 10:25 UTC (12:25 CEST), just four days after launch, this close-up of France's southern coast from Nice airport (lower left) to Menton (upper right) is a subset from the first image from the Sentinel-2A satellite. This false colour image was processed including the instrument's high-resolution infrared spectral channel [35].	32
Figure 20 Landsat 8 Satellite Sensor with a 15m resolution [36]	33
Figure 21 Conceptual rendering of Maxar's WorldView-3 satellite [38]	35
Figure 22 High-resolution satellite measurement by GHGSat of methane emission from an oil and gas facility in New Mexico on September 24, 2022. Taken from [41]	37
Figure 23 They are using VIIRS to pinpoint areas where fuel production is leading to gas flaring, the practice of burning off natural gas that comes to the surface with crude oil, and to estimate the amount of gas released at those sites [45].	39
Figure 24 Detection performance by satellite and team. Total number of measurements listed in brackets. For each satellite, most teams correctly detected most emissions as true positives or true negatives (correctly identified non-emissions). In some cases, e.g. two GHGSat-C2 (GSC2) overpasses, the satellite was not tasked and collected no data. In others, e.g. one SRON retrieval of Landsat 8 (LS8), no retrieval was attempted due to image clipping concerns or excessive cloud cover. No teams produced false positives, in which satellites detected methane when none was released [47]... ..	41
Figure 25 Extreme methane emissions detected in the Permian basin from satellite imaging spectroscopy data. A map with the identified methane plumes is shown in the central panel. Emissions are coded according to their flux rate and to the source of data (GF5-AHSI, GF5; ZY1-AHSI, ZY1; PRISMA, PRS). The small panels (A to J) around the main figure show examples of the detected plumes [48].	48

Figure 26 OMI SO2 SCDs for 16 April 2005 retrieved using(a)the original PCA algorithm and (b) a PCA–NN algorithm, (c) the differences between the two retrievals, and (d) mean SO2 SCDs for 1° latitude bands over relatively clean areas (monthly mean SRR < 3), calculated from (red) the original and (blue) PCA–NN retrievals [49]. 44

Figure 27 Structure of the neural networks used in this study. Green boxes indicate portions of the neural network, orange boxes indicate predictions made by each stage of the neural network. Black lines indicate flow of data into models, and red lines indicate predictions resulting from a model [50]...... 45

Figure 28 Images of plumes detected by the neural network in the Korpjeje oil field, Turkmenistan. Left panels depict methane 376 retrievals, middle panels depict the RGB of the image, and the right panel depicts the mask prediction by the neural network. The predicted emission rates are (top) 7615 and (bottom) 2370 kg hr-1. RGB image courtesy of PRISMA © (Italian Space Agency) [50]...... 46

Figure 29 A schematic overview of the Convolutional Neural Network with a pre-processed 32x32 pixel TROPOMI methane scene (left) as input (Figure 1c). The CNN consists of two convolutional blocks (each with two Convolutional layers followed by a Max-Pooling layer) followed by two Dense (or Fully Connected) layers and an output node. Numerical values show input dimensions, layer dimensions and optimized hyperparameter values [51]. 47

Figure 30 Plumes detected over 10 locations which were inspected with high-resolution instruments. Observations at the same location with different instruments are most often not on the same day [51]. 48

Figure 31 (1) Sentinel-2 Level-1C TOA reflectance input image, (2) the atmospherically corrected Level-2A surface reflectance image, (3) the output scene classification of the Level-1C product. Image taken from [39]...... 51

Figure 32 Monthly average methane from January 2022 including the ocean glint retrieval [40]. 54

Figure 33 overview of the Flaring Monitoring user interface where the sites south of San Antonio has been used as a use case [55]...... 58

Figure 34 A high level overview of the dataset generation 60

Figure 35 An overview of a proposed workflow for developing the pre-processing method based on machine learning.61

Figure 36 Comparison of sentinel-2 data before and after preprocessing: The left image shows an image generated with the hyperspectral data with clouds, while the right image reveals a clear view of the coast after cloud removal. Both images are showing the B4, B3, and B2 bands. 62

Figure 37 Sample images of the Raw Sentinel-2 data with the clouds removed, and with two areas of interest from the VIIRS dataset. On the left, we see a flare site from Kraken Oil and Gas LLC, and on the right, a flare site from Zavanna LLC. 65

Figure 38 Sample image of the processed Sentinel-2 data, with two areas of interest from the VIIRS dataset. On the left, we see a flare site from Kraken Oil and Gas LLC, and on the right, a flare site from Zavanna LLC. 66

Figure 39 Sample image of the NBR index computed from the processed Sentinel-2 data, with two areas of interest from the VIIRS dataset. On the left, we see a flare site from Kraken Oil and Gas LLC, and on the right, a flare site from Zavanna LLC. 67

Figure 40 A representation of the color palette chosen for the Sentinel-5P data, corresponding to a numerical value range of [1750, 1900]...... 68

Figure 41 A visual representation of the integrated Sentinel-2 and Sentinel-5P data, showcasing the combined insights of both platforms. The left image presents the Sentinel-2 multispectral data combined with Sentinel-5P methane concentration data, while the right image displays the Normalized Burn Ratio (NBR) index derived from Sentinel-2 data merged with Sentinel-5P data. Both images shows a flare site for Kraken Oil and Gas LLC..... 68

Figure 42 Detailed overview of the Flaring Monitoring user interface that provides CO2 emission data from US oil and gas sites based VIIRS data 81

LIST OF TABLES



Table 1 Performance Matrix for Satellite Platforms with examples of satellites as input.....	28
Table 2 The landsat-8 specification. Data from [37].....	33
Table 3 WorldView-3 sensor specification. Data from [39]	35
Table 4 Characteristics of the Nine Visible Infrared Imaging Radiometer Suite (VIIRS) Spectral Bands Collecting Data at Night [44].....	38
Table 5 Key characteristics of each participating satellite constellation, from lowest to highest swath width, which is roughly proportional to an instrument’s minimum methane detection limit [47].	41
Table 6 Performance matrix for the evaluated satellite platforms	49
Table 7 Sentinel-2 L2A Product Characteristics	52
Table 8 Sentinel-5 Methane Product Characteristics	55
Table 9 Raw data from the VIIRS instrument.....	58
Table 10 Processed data from the VIIRS instrument.	59
Table 11 An overview of the VIIRS dataset structure, taken from [57].	59



ABBREVIATIONS

ECMWF European Centre for Medium-Range Weather Forecasts

ENVISAT Environmental Monitoring Satellite

GOSAT Greenhousegas Observing SATellite

NASA National Aeronautics and Space Administration

NBR Normalized Burn Ration

netCDF Network Common Data Format

NIR Near InfraRed

NOAA National Oceanic and Atmospheric Administration

OCO OrbitingCarbonObservatory

S5P Sentinel-5Precursor

SWIR Short Wave InfraRed

TROPOMI Tropospheric Monitoring Instrument

UVN Ultraviolet, Visible, Near-Infrared

VIIRS Visible Infrared Imaging Radiometer Suite



INTRODUCTION

In a world where the climate changes rapidly, greenhouse gas emissions significantly drive this transformation. Methane is a potent greenhouse gas with a global warming potential over 80 times greater than carbon dioxide over a 20-year timeframe [1]. Methane emissions from various sources, including oil and gas operations, agriculture, and waste management, have become a growing concern for climate change mitigation efforts. Flaring is a common practice in the oil and gas industry, used to burn off excess gas and prevent dangerous pressure build-up in pipelines. However, flaring also produces significant methane emissions, contributing to climate change.

Methane emissions can originate from various sources, including agriculture, waste management, and fossil fuel extraction, as seen in Figure 3.

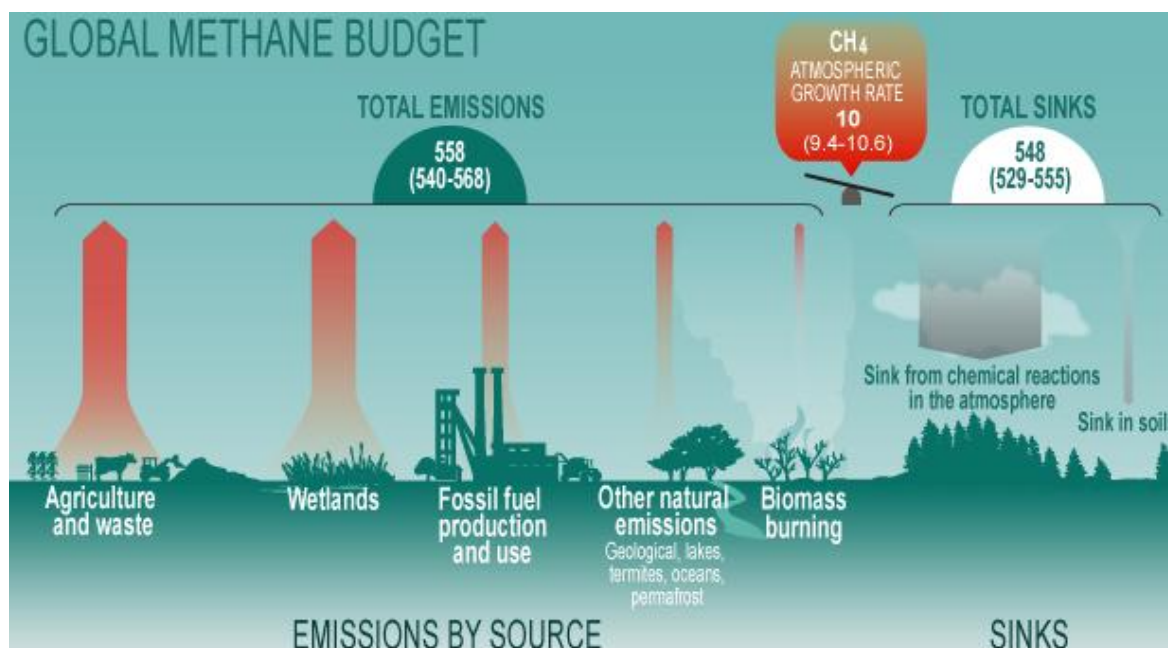


Figure 3 Major sources and sinks for methane. The size of the arrow indicates the relative contribution a source makes to the global total. Methane's lifetime is about nine years before oxidizing agents convert it into carbon dioxide. Image from [2]

The reduction of methane emissions is a top priority, with several key initiatives in place, including the "Reduction Act Methane Emissions Charge," the "Glasgow Agreement," the "European Green Deal," and the "EU's methane strategy." In October 2020, the EU published its methane strategy, which has the potential to significantly enhance its efforts to achieve critical climate objectives, including more ambitious greenhouse gas reduction targets. The strategy seeks to reduce temperature trajectories leading up to 2050, enhance air quality, and reinforce the EU's position as a global leader in combating climate change [3].

The oil and gas industry remains a primary source of methane emissions, with North Sea gas platforms being no exception. Within the oil and gas industry, methane can be released during the production, processing, storage, and transportation of natural gas and petroleum. One common practice in the industry is flaring, which involves the controlled burning of excess gas at production sites. Flaring is intended to reduce the environmental impact of greenhouse gas emissions by converting methane into carbon dioxide, which has a lower global warming potential.

Monitoring methane emissions during flaring operations is essential to quantify the environmental impact of this practice accurately. Methane emissions monitoring can help identify opportunities for reducing emissions, optimize the efficiency of flaring operations, and provide a basis for regulatory compliance. One of the most promising methods for monitoring methane emissions is satellite technology. Satellites equipped with advanced sensors can provide accurate, real-time monitoring of methane emissions over large areas. By analyzing satellite data, researchers and industry professionals can identify potential leaks and assess the effectiveness of mitigation efforts. This information can help reduce methane emissions and improve the efficiency and safety of flaring operations. Monitoring methane emissions during flaring operations is essential for climate change mitigation. Using satellite-based technologies can provide an

effective and efficient means of monitoring methane emissions and identifying opportunities for reducing emissions. Developing and deploying these technologies is critical to ensuring that the oil and gas industry can meet its commitments to reducing greenhouse gas emissions and limiting global warming.

To apply satellite data for implementing a future "Methane Tax," it is essential to improve the accuracy of satellite measurements. This project explores advanced measurement technology as a potentially cost-effective and scalable solution for monitoring methane emissions.

PROBLEM DESCRIPTION

New regulations and taxation schemes have been put forward to reduce methane emissions. For example, in California, legislators have banned new oil wells and mandated a leak detection and response plan for all existing wells within 3200 feet (1 km) from housing, schools, and other special receptors [4]. The bill has been in effect since January 1st, 2023. Thousands of producing oil wells within 3200 feet of these special receptors must develop a leak detection and response plan to be submitted by January 1st, 2025. This situation highlights the added requirements on businesses to monitor their methane releases.

This project focuses on creating a method to independently calculate methane emissions from oil and gas activities to solve this issue. Satellite data, which is helpful remote sensing method for large landmasses, is included in the suggested methodology. In addition to standard monitoring techniques, satellite data offers an independent, non-intrusive, affordable, and continuous monitoring approach.

Based on this, the problem statement for this work is the following

"How can a data-driven approach be developed to enhance the accuracy and quality of methane emission estimation from flaring activities in the Oil and Gas industry, using satellite data from selected platforms to detect and quantify future emissions based on Machine learning more effectively?"

To achieve this, the following objectives and activities are set

- **Theoretical framework and key concepts**
 - What is methane and its effect as a greenhouse gas?
 - What is Flaring?
 - Market drivers
 - What are the key drivers and barriers for methane emission monitoring?
 - Introduction to Satellite technology as a tool for emission monitoring
 - How are methane emissions measured from space?
 - What is Machine Learning?
 - Application of Machine Learning in the Context of methane emission
 - Overview of data quality
 - Summary
 - Define the requirements for a potential monitoring system.
 - Suggest a preliminary model basis for a method.
- **Technical review**
 - «State-of-the-art Analysis»
 - «Present Methane data collection» workflow
 - "Present Flaring emission data collection"
 - Overview of key satellite programs
 - Literature review
 - Identification and review of key research on increasing data quality in satellite emission monitoring.
 - Summary
 - Determine the most suitable method and satellite platform to do further work on
 - Identify a suitable dataset to do further work on
- **Development of the method**

- Review of suitable datasets.
- Process description – «System Description.»
- Development of a «Proof of Concept»
- Summary – presentation of the POC
- **Evaluation of the method**
 - Formulation of test plan
 - Testing of the method
 - Summary – present the results
- **Recommendations and further work**

LIMITATIONS

This thesis aims to develop a data-driven approach for enhancing the accuracy and quality of methane emission estimation from flaring activities in the oil and gas industry using satellite data. While the proposed methodology offers a valuable starting point for addressing this challenge, it is crucial to recognize the limitations associated with the current scope and design of the study.

FOCUS OF THIS WORK

The primary focus of this research is on methane emissions from flaring activities. This limits the scope of the study and may not address other sources of methane emissions in the oil and gas industry, such as fugitive emissions from pipelines, vents, or other equipment. Additionally, the study concentrates on developing a method for improving the accuracy of satellite measurements rather than implementing an end-to-end solution based on machine learning techniques. This limitation may require further refining and expanding the methodology to include advanced predictive models and algorithms.

SATELLITE DATA SELECTION

The thesis limits its analysis to several satellite platforms, including Sentinel-5, Sentinel-2, Landsat-8, WorldView 8, GHGSAT-METHANESAT, and the Visible Infrared Imaging Radiometer Suite (VIIRS). While these satellites provide valuable data for estimating methane emissions, excluding other platforms may limit the comprehensiveness and accuracy of the analysis. Additionally, each satellite has inherent limitations, such as spatial and temporal resolution, coverage, and data quality, which may impact the overall effectiveness of the proposed methodology.

EVALUATION CONSTRAINTS

As the focus of this thesis is on creating a method for generating a dataset, the project limits itself to proposing a test plan for evaluating the accuracy and quality of the generated dataset and methodology, rather than performing an in-depth evaluation. A comprehensive evaluation, including ground truth data collection, comparison with existing estimation techniques, data quality assessment, and analysis of potential implications for the oil and gas industry, would provide a more robust understanding of the strengths and weaknesses of the proposed method.

In summing up, while the methodology outlined in this thesis offers a promising strategy for estimating methane emissions from flaring activities using satellite data, it's crucial to consider the limitations related to the scope of this work, the selection of satellite data, and evaluation constraints. Addressing these limitations in future iterations of this methodology would help enhance its refinement and broaden its potential applications within the oil and gas industry.

THEORY AND KEY CONCEPTS

The purpose of this section is to provide an overview of key concepts related to developing a machine learning-based method for improving emission estimation from flaring activities in the Oil and Gas industry using satellite data. The chapter covers methane's environmental impact, flaring operations, market drivers for monitoring emissions, satellite technology and its applications, supervised machine learning methods, data quality, model assessment, and satellite data quality considerations such as resolution, accuracy, spectral range, and signal-to-noise ratio. This foundational knowledge paves the way for devising an effective solution for enhanced methane emission estimation from flaring activities.

METHANE AND ITS ENVIRONMENTAL IMPACT

Methane (CH_4) is a colorless, odorless gas that is the primary component of natural gas and is the second most important greenhouse gas contributor to climate change, following carbon dioxide. It is a potent greenhouse gas with a global warming potential that is 23 times greater than carbon dioxide (CO_2) over a 100-year time frame and is 84 times more potent on a 20-year timescale. Therefore, methane emissions are highly pertinent to the climate objectives set for 2050, not least those established by the European Union. Moreover, methane is a potent local air pollution and contributor to ozone formation, which causes serious health problems [5] [6].

Atmospheric methane has many natural and anthropogenic emission sources. These sources can be split into three main types: microbial methane, which is emitted during the decomposition of organic matter; pyrogenic methane, which is formed during incomplete combustion of biomass; and thermogenic methane, which is released from fossil fuels during their extraction, refinement, and use. Globally, anthropogenic sources account for approximately 60 % of total methane emissions. The most significant sources of methane are agriculture and waste management (approximately 35 % of total emissions) and fossil fuel production and use (approximately 20 % of total emissions).

Figure 2 presents a comprehensive illustration of methane emissions from fossil fuel exploitation [3]. This figure provides valuable insights into the various sources and contributing factors that lead to these emissions, enabling a better understanding of their overall environmental impact.

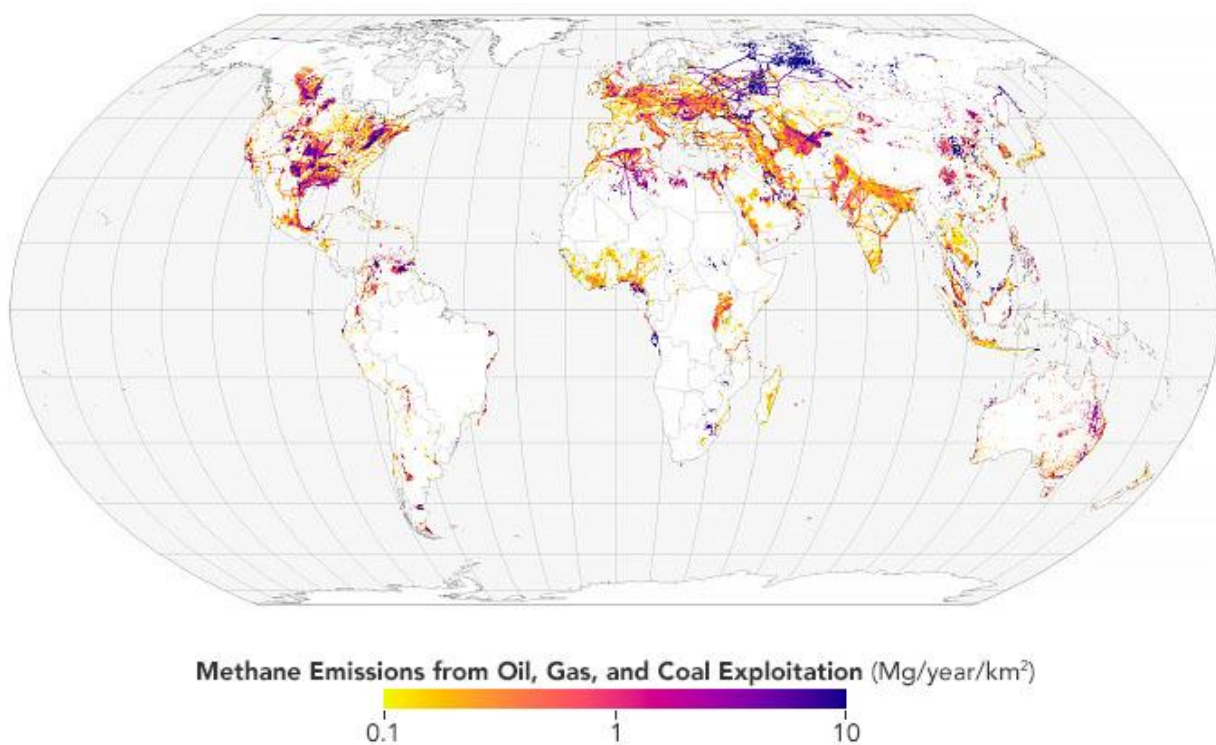


Figure 4 An overview of methane emissions from fossil fuel exploitation [7] .

WHAT IS FLARING AND THE EMISSION FROM FLARING OPERATIONS?

Flaring is the process by which natural gas is burned off in a controlled manner when extracting oil. Otherwise, the natural gas can burn in an uncontrolled way and be very dangerous. Usually, natural gas is captured, but when this is impossible, it is flared. Flaring may be required for safety reasons as it reduces the risk of gas ignition to facilities or eliminates product that is not fit for use [8].

In some cases, it is economically and technically feasible to capture and utilize associated gas. However, a country's laws and regulations might make it difficult for, or even forbid, companies from selling associated gas. Figure 5 shows a typical flare system.

In 2021, roughly 144 billion cubic meters of gas were burned by thousands of gas flares at oil production sites worldwide. With a standard associated gas composition, a 98% flare combustion efficiency, and a methane Global Warming Potential of 25, each cubic meter of flared gas generates around 2.8 kilograms of CO₂ equivalent emissions. This amounts to over 400 million tons of CO₂ equivalent emissions annually. Methane emissions from flare combustion inefficiency significantly impact global warming, especially in the short to medium term, as methane is over 80 times more potent than carbon dioxide as a warming gas in a 20-year timeframe, according to the Intergovernmental Panel on Climate Change. Consequently, this increases annual CO₂ equivalent emissions by nearly 100 million tons.

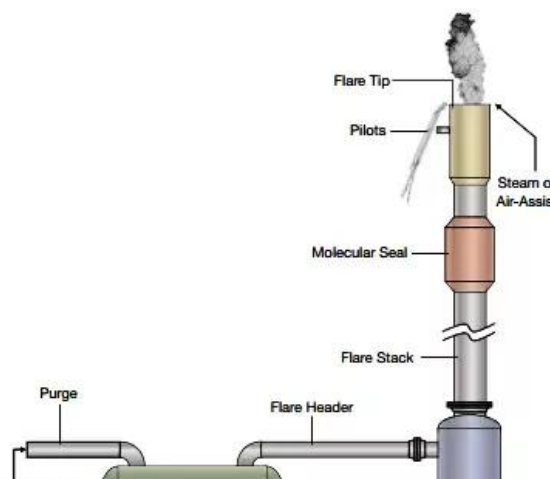


Figure 5 Overview of typical gas Flare system

In Norway, there has been a carbon tax on flaring since the 1990s, significantly reducing the amount of flaring in the Norwegian oil and gas sector. This has mainly been enabled with flare meters based on ultrasonic meters that measure the amount of gas being burnt per hour, which is then taxed [9].

Various initiatives, like the "Zero Routine Flaring by 2030" launched in 2015, aim to reduce flaring activities. These initiatives involve commitments from governments and oil companies to end routine flaring by 2030 [8]. However, despite these efforts, significant emissions from flaring still exist. Satellite data could be crucial in monitoring flaring activities and supporting these initiatives.

While, in theory, the flaring of gas from oil wells and installations is supposed to combust methane and release only carbon dioxide, the reality is often different. Flares are frequently inefficient and do not fully combust methane. It has been traditionally assumed that flaring is 98% efficient in converting methane to carbon dioxide [8]. However, recent research by a team from the University of Michigan, Stanford, the Environmental Defense Fund, and Scientific Aviation challenges this assumption. Their findings suggest that flaring efficiency could be as low as 91.1% [9]. This discrepancy has significant implications for estimating methane emissions and the effectiveness of current mitigation strategies. It means that a notable portion of methane is released into the atmosphere instead of being converted to carbon dioxide. Unlit or malfunctioning flares also contribute to the problem. Additionally, flaring is often underreported, emphasizing the need for improved monitoring capabilities for methane emissions.

WHAT ARE THE MARKET DRIVERS FOR MONITORING METHANE EMISSIONS?

The overall question for the oil and gas industry is perhaps, *“How can the oil and gas industry—and the regulators who oversee it—best detect and address methane emissions to protect the environment and the climate in particular?”*

The key market drivers for improved monitoring of methane emission can be divided into the following sections

- Drivers for
 - policy development
 - oil and gas companies
 - Investors
 - Public

Companies responsible for production, processing, transportation, and distribution in the oil and gas industry often self-report their methane emissions, even when participating in voluntary emissions-reduction programs. This self-reporting makes it challenging for external parties to verify the accuracy of the reported data. The motivations behind these efforts vary among companies, with limited incentives ranging from good business practices to maintain the industry's social license to operate. Some companies also recognize the value of minimizing methane waste, as it represents a valuable resource.

However, when official emissions figures are based on engineering estimates that seem to underreport actual emissions levels systematically, there is minimal external pressure for transparency, accuracy, or addressing actual emissions [10]. The recent decline in global demand for oil and natural gas has made it even more challenging for companies that have not been addressing methane emissions to start doing so now. The industry faces a significant loss of trust from the public, and as awareness of actual emissions levels grows, the environmental credibility of natural gas is being scrutinized more skeptically.

The emergence of progressively improving satellite-based methane emissions information will subject the oil and gas industry to increased scrutiny from stakeholders with more accurate and timely data. This increased attention will benefit the climate, as it will cover the industry's entire value chain, from exploration and production to local distribution. It will also apply to non-operated assets and joint ventures, including jurisdictions with low transparency and environmental standards. This is likely to intensify controversies over new developments and ongoing operations. Enhanced satellite-based emissions data may facilitate comparisons between companies in different jurisdictions, making it easier to assess their performance, regardless of whether they are publicly traded "majors," smaller independent companies, or state-owned entities [10].

IMPLICATIONS FOR INVESTORS

According to an analysis by the United States Federal Reserve System, the financial and investment community is becoming more concerned with understanding and addressing vulnerabilities related to climate change causes and impacts [11]. Currently, investors interested in basing their decisions on environmental risk and opportunity assessments often lack accurate data to evaluate companies' performance. They typically rely on environmental, social, and governance (ESG) ratings for potential investments. However, today's ESG ratings can cause more confusion and complexity than providing clear and definitive insights. With over 600 ESG ratings and rankings available, the landscape is fragmented. Various rating schemes often use non-transparent methodologies or "black boxes" without standardized scoring approaches. This issue is exacerbated by the fact that ESG rating firms must rely on unreliable public emissions data. Consequently, a company's rating can significantly vary between different rating agencies [11].

IMPLICATIONS FOR PUBLIC POLICY

The public policy impacts of a more accurate and timely understanding of methane emissions could prove especially significant.

This better understanding may well reverberate through national and international structures whose mission it is to maintain inventories of greenhouse gas emissions and craft policies to respond adequately to them—including national-level environment agencies, scientific entities like the Intergovernmental Panel on Climate Change, and the parties to the

UN Framework Convention on Climate Change. If it turns out that global emissions of oil-and-gas-related methane are significantly greater than previously understood, national and international priorities for climate change mitigation may need to be reordered.

The availability of improved methane emissions data could encourage countries to demand that their companies detect and respond to methane emissions. The new data could drive efforts to introduce strict regulation of gas production or imports or to accelerate the abandonment of fossil fuel usage. The new data could also facilitate emissions pricing schemes. The corollary to this thought is that, because of methane's potency and short atmospheric lifetime, addressing methane emissions with a much greater sense of urgency may help to foster badly needed climate progress. Reducing and eliminating methane emissions, after all, should enable more accessible and faster climate mitigation than alternative options, such as significantly reducing emissions from heavy industry or replacing internal combustion engines with electric vehicles.

Other impacts in public policy will arise regarding current legal and regulatory structures. Policymakers, regulators, the general public, and environmental advocacy groups that seek to speak on the public's behalf will have a much better understanding than was historically the case of how much methane is being emitted, where, by whom, and for how long. This information may be employed to identify ineffective regulatory oversight or to "name and shame" individual companies whose operations result in methane emissions. The data may be introduced into facility-level approvals as well as long-term policy planning. In this sense, the new information may assist regulators in enforcing emissions standards, penalizing laggards, or, conversely, creating greater incentives for more robust environmental stewardship.

POLICY DEVELOPMENT

Lastly, the new availability of methane emissions data may be used to exert influence on policymakers and regulators themselves, as seen by the "Inflation Reduction Act Methane Emissions Charge:" in the US and "The EU Methane Strategy" as part of the "European Green Deal."

The Inflation Reduction Act Methane Emissions Charge imposes a first-ever direct "charge" on methane emissions. The charge is based on an unworkable comparison between the weight of methane emitted and the volume of the natural gas stream sent to sale. It states that *"The charge starts at \$900 per metric ton of methane, increasing to \$1,500 after two years, which equates to \$36 and \$60 per metric ton of carbon dioxide equivalent, respectively. This charge is the first time the federal government has directly imposed a charge, fee, or tax on GHG emissions"* [12]. Cost Estimate Analysis shows that this will generate between \$ 500 to 1800 million in tax revenue from the US industry [13].

These drivers are fueling the development of improved measuring methods, such as sensors and methane emission monitoring via drones but also satellites.

OVERVIEW OF SATELLITE TECHNOLOGY

Satellites are objects orbiting larger objects than themselves. It is divided into two groups: natural satellites and artificial satellites. Satellites developed and placed in orbits of planets are called artificial satellites. Sputnik 1, the first artificial satellite in history, was launched by the USSR from Earth in 1957.

According to [14], there are six types of orbits for satellites:

1. Low Earth Orbit (LEO)
2. Geosynchronous Earth Orbit (GEO)
3. Medium Earth Orbit (MEO)
4. Polar orbit and Syn-synchronous orbit (SSO)
5. Transfer orbits and geostationary transfer orbits (GTO)
6. Lagrange points (L-points)

And a majority of satellites, about 55%, are in Low Earth Orbit [15].



Figure 6 Example of a satellite in orbit around the Earth. Taken from [16]

The equipment on a satellite is called a "Payload," The term 'Payload' was originally a seafaring term for revenue-producing cargo on a ship [17]. In space terms, it refers to those elements of the spacecraft dedicated explicitly to producing mission data and then relaying that data back to Earth.

Today's technology provides people with many opportunities in many areas. One of these areas is satellite technologies, that are increasingly developing today. Especially in recent years, satellite technologies have provided fast, economical, and high-accuracy information on the Earth. Earth can be imaged in different time intervals and different spatial and spectral resolutions with remote sensing technologies, and the satellite images obtained can be used in many studies. Thanks to geographic information systems, the integration of satellite images and different location data can be made. In this way, analysis, interpretation, and solution suggestions for many location-based problems can be provided.

THE APPLICATION OF SATELLITES IN METHANE EMISSION MONITORING

A satellite payload is the equipment or instruments that a satellite carries into space to perform a specific function or mission. The payload is part of the satellite that is responsible for collecting data, transmitting signals, or performing any other task that the satellite is designed to do.

Satellite payloads, as seen in Figure 7, can include a wide range of equipment, such as cameras, sensors, communication systems, scientific instruments, and more. The specific payload that a satellite carries depends on its mission, whether it is to observe the Earth's surface, study the atmosphere, monitor weather patterns, provide global positioning services, or any number of other functions. The payload is typically mounted onto the satellite's bus, which is the framework that provides the satellite with power, propulsion, and other essential systems. Once in orbit, the payload operates independently of the bus, collecting data or performing other functions and transmitting the information back to Earth.

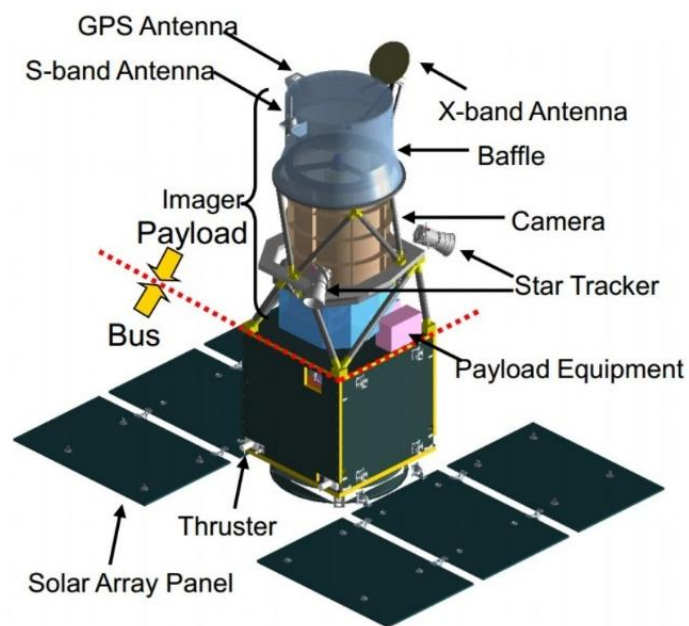


Figure 7 Overview of a typical satellite Payload. Image taken from [18].

A prototypical sensor payload within a satellite embodies a collection of instruments or apparatuses contrived to detect and quantify diverse physical attributes or phenomena. These sensors can be categorized as passive or active, depending on the methodology of data collection. The specific types of sensors included in the payload rely on the satellite's mission and objectives.

The following sensors may comprise a typical sensor payload:

- Optical sensors
- Synthetic Aperture Radar (SAR) Sensors
- Multispectral and Hyperspectral Sensors
- Radiometers
- Lidar Sensors
- Magnetometers
- GPS Receivers

These payloads enable remote sensing and can be utilized for scientific research, natural resource management, disaster monitoring, and a myriad of other applications.

SHORT-WAVE INFRARED

Methane has absorption lines in the infrared, as shown in Figure 8, which can be utilized to detect and quantify methane using spectroscopy.

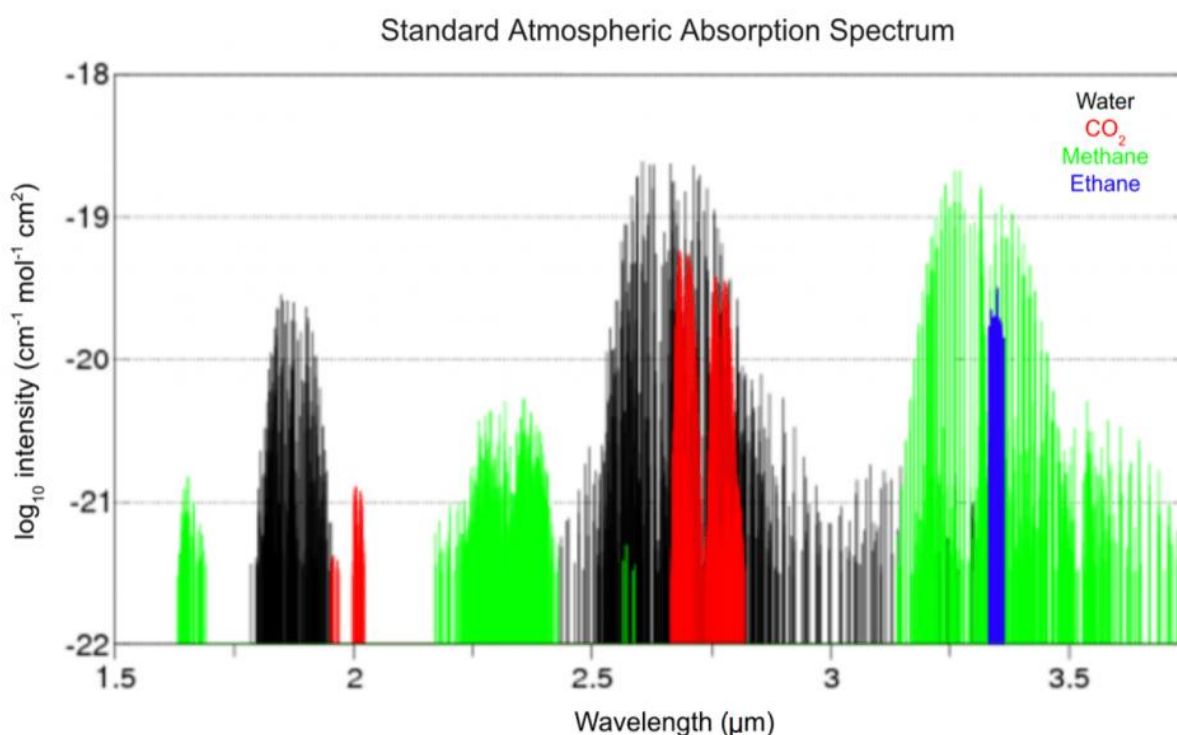


Figure 8 Absorption spectrum of methane (green), water (black), carbon dioxide (red), and ethane (blue) in the infrared region of the electromagnetic spectrum [19].

Methane's absorption lines are uniquely identifiable, setting them apart from other gaseous species. Consequently, spectroscopic analysis can effectively distinguish methane from other potentially interfering elements like water, carbon dioxide, and ethane, as demonstrated in Figure 8.

When sunlight that has been reflected travels through a rogue methane plume in an oil and gas field, specific infrared wavelengths are absorbed by the gas molecules. This sunlight, now reflected from the ground, can be captured via aerial surveillance. The execution of a spectroscopic analysis on this reflected sunlight allows for the detection of excessive methane originating from flaring activities or leaking infrastructure.

HYPERSPECTRAL SENSING

Hyperspectral sensing is a remote sensing technique that captures high-resolution spectral data across a wide range of wavelengths. In the context of methane monitoring, hyperspectral sensing is valuable for detecting and quantifying methane emissions from various sources, including natural gas leaks, landfill sites, and agricultural operations.

Methane gas has specific absorption patterns at certain wavelengths, enabling the use of hyperspectral imaging to identify and measure methane concentrations. Hyperspectral sensors can capture detailed images of methane plumes, which can then be analyzed to estimate the quantity of methane released. By comparing images taken at different times, changes in methane emissions can be monitored, facilitating early detection and swift response to leaks or other methane emission sources.

IMAGE ENHANCEMENT METHODS FOR HYPERSPECTRAL DATA

NORMALIZED BURN RATIO

The Normalized Burn Ratio (NBR) index is a widely used remote sensing technique for assessing the severity of burned areas and monitoring post-fire vegetation recovery. The NBR index is derived from satellite data, particularly from multispectral sensors like those on the Sentinel-2 satellite.

The NBR index is calculated using the near-infrared (NIR) and shortwave infrared (SWIR) bands of the satellite imagery. The formula for NBR is as follows:

$$NBR = \frac{NIR - SWIR}{NIR + SWIR} \quad (1)$$

The rationale behind using these two bands lies in their distinct responses to vegetation and burned areas. Healthy vegetation has a high reflectance in the NIR region due to its cell structure and a low reflectance in the SWIR region because of the water content in the leaves. In contrast, burned areas and regions affected by gas flaring exhibit low reflectance in the NIR region and higher reflectance in the SWIR region, mainly due to the loss of vegetation and the presence of charred materials.

The NBR index generates values ranging from -1 to 1. Higher positive values indicate healthy vegetation, while lower negative values signify burned areas or areas with little to no vegetation. The NBR index is particularly useful for estimating the extent and severity of fire damage and for tracking the progress of vegetation recovery over time.

In the context of the dataset generated from this project, the NBR index can provide valuable information on the impact of flaring events on the surrounding environment. Analyzing changes in the NBR index over time can help the user of the dataset to infer the effects of flaring on vegetation and land cover, which in turn can provide insights into the relationship between methane emissions and gas flaring activities. Monitoring these changes can contribute to a better understanding of the environmental consequences of flaring and potentially help models map the amount of environmental damage around a flare site to the amount being flared.

SPATIAL BINNING OF SATELLITE DATA

Spatial binning is a widely used process in satellite data processing and analysis, which involves the aggregation of data points within a specific spatial domain, often defined by a grid of equally-sized cells or bins. This process is critical for transforming Level 2 (swath or granule) data into Level 3 (gridded) data, which is easier to analyze, visualize, and integrate with other spatial datasets.

Spatial binning is fundamentally a process of spatial aggregation, where multiple data points within a given spatial extent (i.e., the bin) are summarized into a single value that represents the bin as a whole. This process is underpinned by the theories of spatial statistics, which deals with the analysis and interpretation of spatially referenced data.

The most basic form of spatial binning involves averaging the values within each bin, which can be expressed mathematically as:

$$V_{bin} = \frac{1}{N} * \sum V_i \quad (2)$$

Where V_{bin} is the binned value, N is the number of points in the bin, and V_i is the value of the i -th point in the bin.

However, more complex aggregation operations can also be applied, depending on the nature of the data and the specific requirements of the analysis. For example, one might choose to take the median, mode, maximum, or minimum value within each bin, or to calculate a measure of variance or dispersion.

The size and extent of the bins are critical parameters in the binning process. The size of the bins determines the spatial resolution of the resulting Level 3 dataset, while the extent of the bins (i.e., the range of latitudes and longitudes covered by the grid) determines the geographical coverage of the dataset.

The bin size is a trade-off between spatial resolution and noise reduction. Larger bins provide greater noise reduction, as the averaging process smooths out random variation in the data, but at the cost of reducing the spatial resolution of the dataset. Conversely, smaller bins retain more spatial detail but may result in a noisier dataset.

The grid extent, on the other hand, should be chosen to match the geographical area of interest for the analysis. Data points outside this area can be excluded from the binning process, which can significantly reduce the amount of data that needs to be processed and stored.

INTRODUCTION TO MACHINE LEARNING

In this thesis, we will explore the development of a dataset suitable for machine learning-based method aimed at enhancing the accuracy and data quality of emission estimations from flaring activities in the Oil and Gas industry using satellite observations. The rationale behind employing machine learning techniques is to create robust estimators that can effectively determine emissions without the need for intricate knowledge about the specific properties of methane gas.

WHAT IS MACHINE LEARNING?

Machine learning, a subcategory of artificial intelligence (AI), centers on creating algorithms and computational models that empower computers to learn from data and use it to predict or decide. It encompasses the procedure of automatically detecting patterns, extracting significant findings, and rendering data-influenced decisions without requiring explicit programming for every individual task.

There are three primary types of machine learning.

1. Supervised learning
2. Unsupervised learning
3. Reinforcement learning

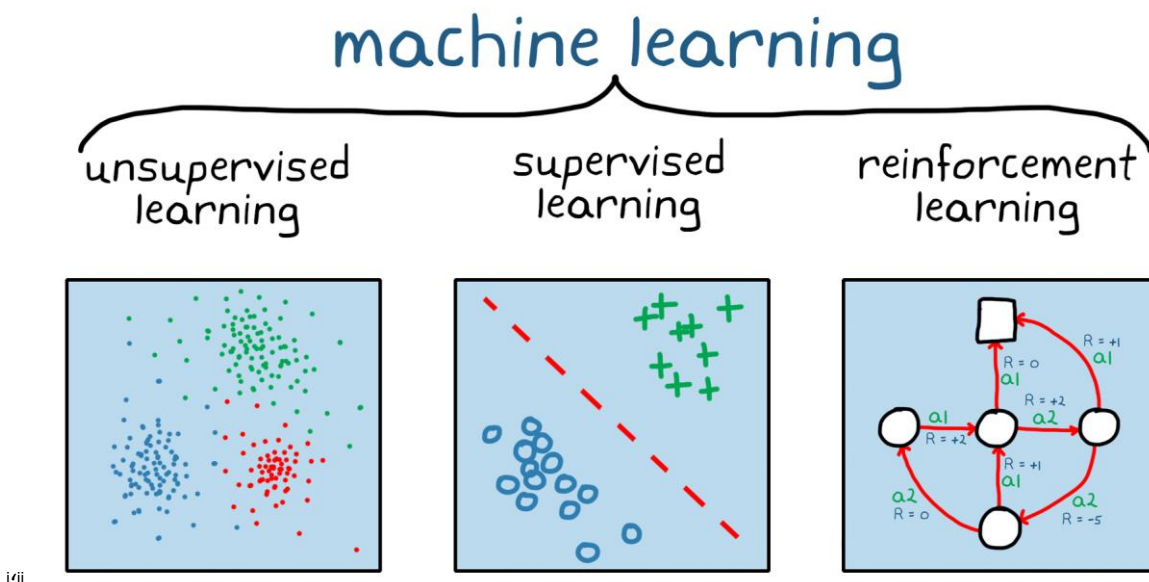


Figure 9 Three primary types of machine learning. Image taken from [20].

SUPERVISED LEARNING

Supervised learning is a type of machine learning where the algorithm learns to make predictions or decisions based on labeled training data. The labeled data consists of input features and corresponding output labels, and the algorithm tries to learn the mapping between the input and output by minimizing a loss function.

The underlying mathematics of supervised learning involves the use of statistical models and optimization techniques. The goal is to find a function that maps the input features to the output labels. The function can be represented by a mathematical model such as a linear regression model, a decision tree, a neural network, or a support vector machine (SVM). Usually, this is a cost function, a cost function is a mathematical function used to measure the difference between the predicted output and the actual output in a machine-learning model. The cost function helps to evaluate how well the model is performing and is used as a basis for the optimization algorithm to update the model parameters during the training process.

In supervised learning, the cost function is typically defined as the average difference between the predicted output and the actual output over the training dataset. The goal is to minimize the cost function by adjusting the model parameters, such as the weights and biases in a neural network.

There are different types of cost functions used in machine learning depending on the type of problem being solved. For example, in regression problems, the mean squared error (MSE) cost function is commonly used, as seen in Equation 3

$$MSE = \frac{1}{N} \sum_{i=1}^N (\text{actual value}(y) - \text{predicted value}(x))^2 \quad (3)$$

MSE measures the average squared difference between the predicted and actual output. In classification problems, the cross-entropy cost function is commonly used, which measures the difference between the predicted probability distribution and the actual probability distribution.

The choice of cost function is important because it affects the performance of the model and the convergence of the optimization algorithm. A good cost function should be differentiable and convex and should have a unique global minimum. The optimization algorithm works by iteratively adjusting the model parameters in the direction of the negative gradient of the cost function, so the cost function should be differentiable to enable gradient-based optimization techniques to be used.

To train a machine learning model, one can use a training dataset that consists of input features and corresponding output labels. The algorithm tries to learn the mapping between the input and output by adjusting the model parameters. This is done by minimizing a loss function, which measures the difference between the predicted output and the actual output.

The most common optimization technique used in supervised learning is gradient descent, which iteratively adjusts the model parameters to minimize the loss function. The gradient descent algorithm computes the gradient of the loss function with respect to the model parameters and updates the parameters in the direction of the negative gradient.

Once the model is trained, it can be used to make predictions on new data that the model has not seen before. The model takes the input features as input and outputs the predicted output label based on the learned mapping.

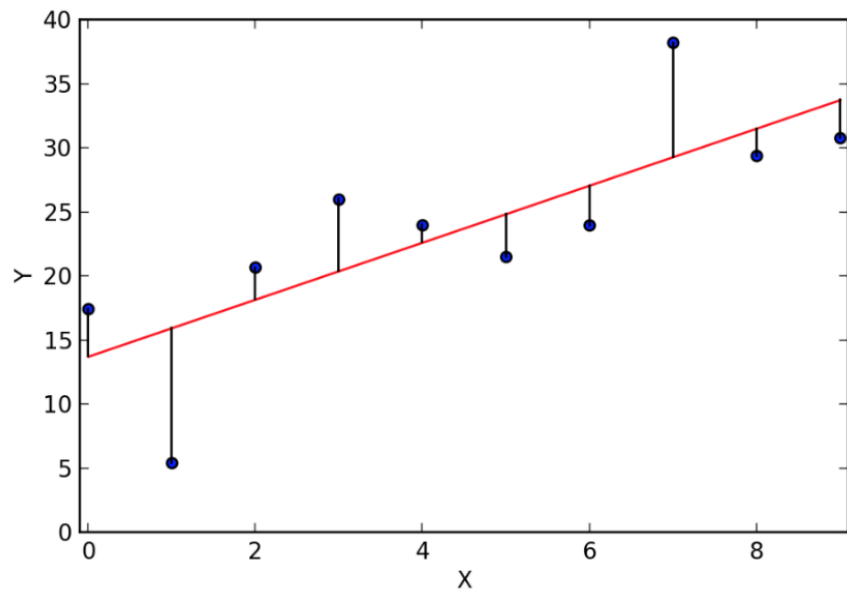


Figure 10 A cost function - In the above image the farther the points is from the straight red line higher the error in predicted value w.r.t ground truth value. Figure taken from [70].

UNSUPERVISED LEARNING

Unsupervised learning is a type of machine learning where the algorithm learns to identify patterns or structures in the data without relying on labeled training data. In contrast to supervised learning, unsupervised learning algorithms work with input features only, without corresponding output labels, and the objective is to discover the underlying structure in data by minimizing a specific objective function [21].

Unsupervised learning's core mathematics revolves around employing statistical models and optimization methods to uncover concealed structures or relationships within the data. Typical unsupervised learning approaches encompass clustering techniques like k-means and hierarchical clustering, as well as dimensionality reduction methods, such as principal component analysis (PCA). These techniques frequently depend on either a similarity or distance metric for grouping or dimensionality reduction purposes [21].

In unsupervised learning, the objective function is designed to measure the quality of the discovered structure, such as the compactness of clusters or the preservation of local distances, in dimensionality reduction. The goal is to minimize the objective function by adjusting the algorithm parameters, like the cluster centroids in k-means or the low-dimensional representations in PCA [21].

In the realm of clustering problems, the predominant objective function involves minimizing the within-cluster sum of squares (WCSS), alternately known as inertia or distortion. This objective function constitutes the foundation of the widely recognized k-means clustering algorithm. The primary goal of this algorithm is to assign data points to clusters in such a way that the aggregate of squared distances between each data point and its corresponding cluster centroid is minimized.

$$WCSS = \sum_{j=1}^k \sum_{x_i \in C_j} \|x_i - c_j\|^2 \quad (4)$$

In this equation:

- k is the number of clusters
- x_i represents a data point in the dataset
- C_j represents the set of data points assigned to cluster j
- c_j represents the centroid of cluster j
- $\|x_i - c_j\|^2$ is the squared Euclidean distance between x_i and c_j

By minimizing the WCSS, the algorithm achieves more compact and homogeneous clusters, wherein data points within each cluster are near their respective centroids. However, it is crucial to acknowledge that the k-means algorithm is sensitive to the initial positioning of cluster centroids and may converge to local optima. To circumvent this issue, the algorithm is typically executed multiple times with varying initializations, and the solution with the lowest WCSS is chosen as the result.

WCSS function quantifies the sum of squared distances between data points and their corresponding cluster centroids. Selecting an appropriate objective function is crucial, as it influences the performance of the algorithm and the convergence of the optimization process. A good objective function should ideally be differentiable, have a meaningful interpretation, and exhibit desirable properties like local or global optima. The optimization algorithm operates by iteratively adjusting the algorithm parameters in the direction that minimizes the objective function, so differentiability is often an important requirement to facilitate gradient-based optimization techniques.

Unsupervised learning algorithms are trained using a dataset consisting solely of input features without corresponding output labels. The algorithm attempts to discover hidden structures in the data by optimizing the algorithm parameters to minimize the objective function. This is achieved using various optimization techniques, such as gradient descent, expectation maximization, or other iterative refinement methods.

Once the unsupervised learning algorithm has been trained, it can be utilized to analyze new, unseen data and extract meaningful insights, such as grouping similar data points together or reducing the dimensionality of the data for visualization or further processing. By leveraging the learned structure, unsupervised learning algorithms can provide valuable information about the underlying patterns and relationships within the data.

REINFORCEMENT LEARNING

Reinforcement learning is a type of machine learning where an agent learns to take actions in an environment to maximize a cumulative reward signal. Unlike supervised learning, the agent is not given explicit examples of the correct action to take but rather learns through trial and error.

The basic idea behind reinforcement learning is that the agent interacts with the environment by taking actions, receiving a reward signal from the environment, and updating its policy to improve its future actions. The policy is a function that maps states of the environment to actions, and the goal of the agent is to learn a policy that maximizes the expected cumulative reward over time.

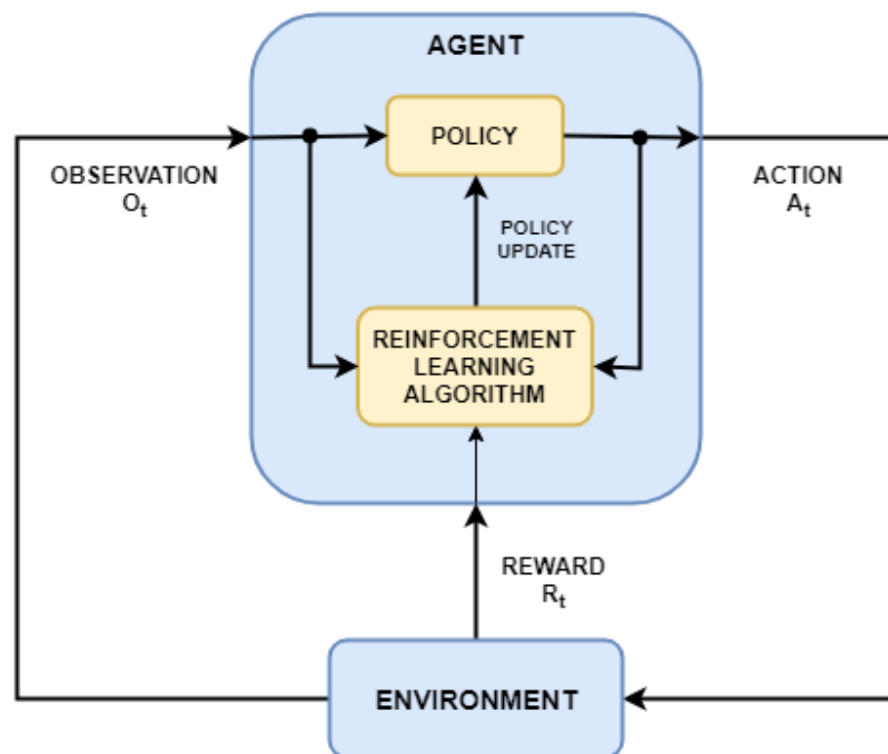


Figure 11 General representation of a reinforcement learning scenario [20].

The underlying mathematics of reinforcement learning involves the use of Markov decision processes (MDPs) and Bellman equations. MDPs are a mathematical framework for modeling decision-making problems where the outcome of an action depends on the current state of the environment. Bellman equations are a set of recursive equations that relate the value of a state or state-action pair to the values of its neighboring states or state-action pairs. These equations form the basis for many reinforcement learning algorithms, including Q-learning and policy gradient methods.

In Q-learning, the agent learns an action-value function that estimates the expected cumulative reward of taking a particular action in a particular state. The action-value function is updated using the Bellman equation, which relates the value of a state-action pair to the values of its neighboring state-action pairs. The policy is then derived from the action-value function by selecting the action with the highest expected reward.

In policy gradient methods, the agent learns a policy directly by optimizing a parameterized policy function to maximize the expected cumulative reward. The gradient of the expected cumulative reward with respect to the policy parameters is computed using the policy gradient theorem, and the parameters are updated using stochastic gradient descent.

METHODS OF SUPERVISED MACHINE LEARNING

In this section, we present the main supervised machine-learning methods suitable for estimating methane emissions from hyperspectral satellite images. These techniques offer distinct advantages, such as accuracy, efficiency, and adaptability, enabling the efficient extraction of information from complex datasets and providing accurate estimations of gas emissions. By employing these methods, researchers could effectively monitor methane emissions on a global scale, addressing the challenges posed by remote sensing data and contributing to our understanding of climate change.

REGRESSION

Regression is a statistical technique that models the relationship between a dependent variable (target) and one or more independent variables (predictors or features). The objective of regression is to forecast the dependent variable's value based on the independent variables [21].

The key underlying mathematics in regression is the estimation of the coefficients of a mathematical equation that describes the relationship between the independent and dependent variables. In linear regression, for example, the equation is a straight line:

$$y = b_0 + b_1x_1 + \dots + b_mx_m = \sum_{i=0}^m b_ix_i = b^t x \quad (5)$$

where y is the dependent variable, x_1, x_2, \dots, x_n are the independent variables, b_0 is the intercept or bias term, and b_1, b_2, \dots, b_n are the coefficients or weights associated with each independent variable. The goal of the algorithm is to find the values of the coefficients that minimize the difference between the predicted values and the actual values in the training data.

To estimate the coefficients, regression algorithms use a variety of optimization techniques, such as least squares or gradient descent. The optimization process involves minimizing a cost function, as previously described, that measures the difference between the predicted values and the actual values. The cost function can be defined in different ways, depending on the type of regression and the specific problem being solved. The most common cost function in linear regression is the mean squared error (MSE) which is shown by equation (3).

Regression models are used in a variety of applications, such as predicting housing prices based on features like location, size, and number of bedrooms; forecasting stock prices based on historical data; and estimating the probability of a patient developing a certain disease based on their medical history and lifestyle factors.

Regression encompasses several algorithms, such as linear, logistic, polynomial, ridge regression, and more [21]. Each algorithm has its strengths and limitations, and the selection depends on the dataset's characteristics and the problem being solved.

NEURAL NETWORK

Neural networks, also known as artificial neural networks (ANNs) or simulated neural networks (SNNs), are a subset of machine learning and are at the heart of deep learning algorithms [22]. A neural network is a type of machine learning model that is inspired by the structure and function of the human brain. It is a collection of interconnected nodes (also called neurons) that are organized in layers, and these layers allow the network to learn from input data and make predictions or decisions based on that data. As seen in Figure 12 below, the neural network consists of three types of layers: input layer, hidden layer, and output layer. The input layer receives the input data, and the output layer produces the final output of the network. The hidden layer(s) perform computations on the input data and learn to represent the underlying patterns in the data.

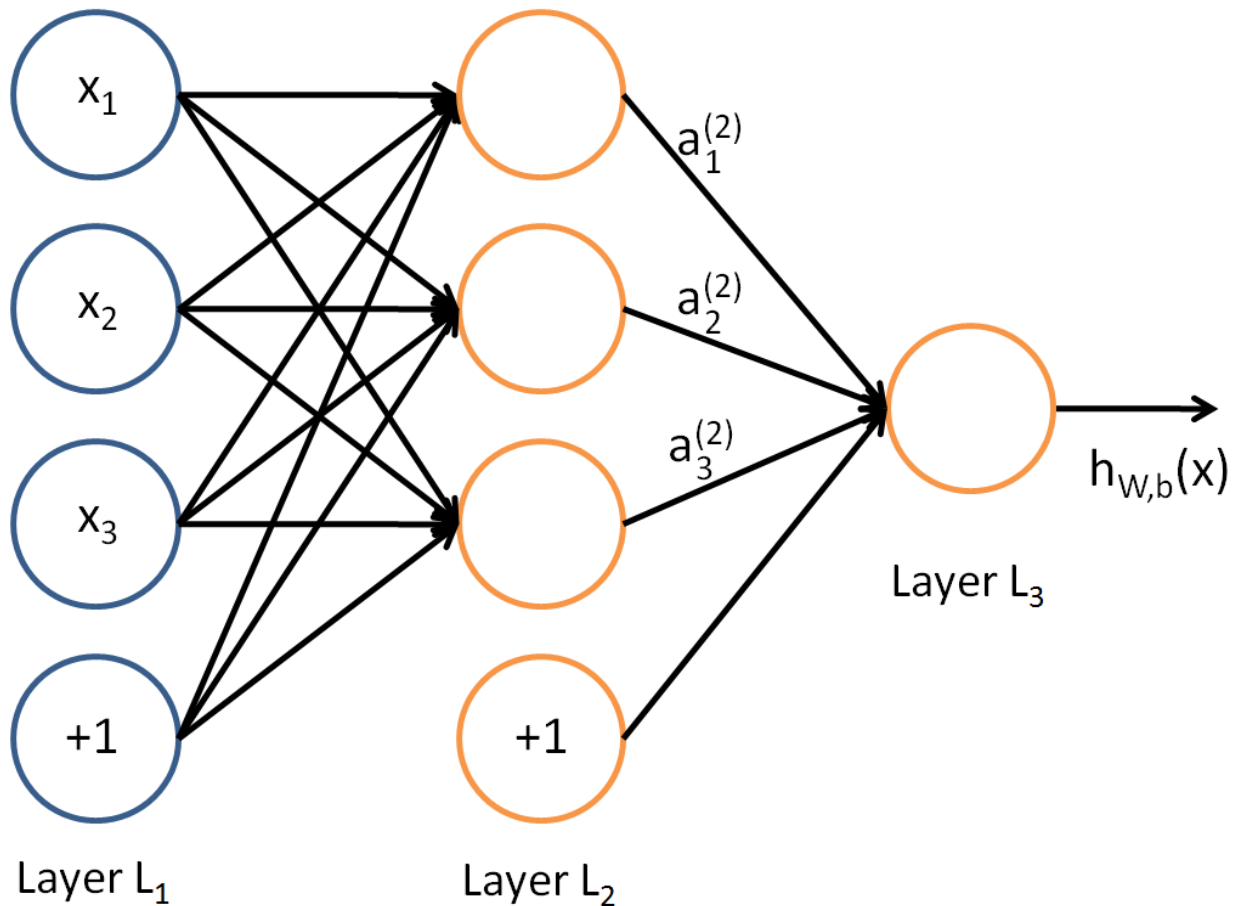


Figure 12 Conceptual overview of a neural network [23].

Figure 12 shows a conceptual overview of a neural network, with each neuron shown as a circle. The neural network consists of multiple layers, denoted by L_n , where n denotes the layer of the network. Each neuron in the network is denoted by x_i , where i denotes the index of the neuron. The circles denoted with a +1 are bias units and represent the intercept term in the context of neural networks. They are added to every layer before the output layer without any connections to any previous layer.

Each neuron in a neural network is connected to other neurons through a set of weighted connections or can be viewed as each individual node as its own linear regression model, composed of input data, weights, a bias (or threshold), and an output. The following illustration shows a single “neuron.”

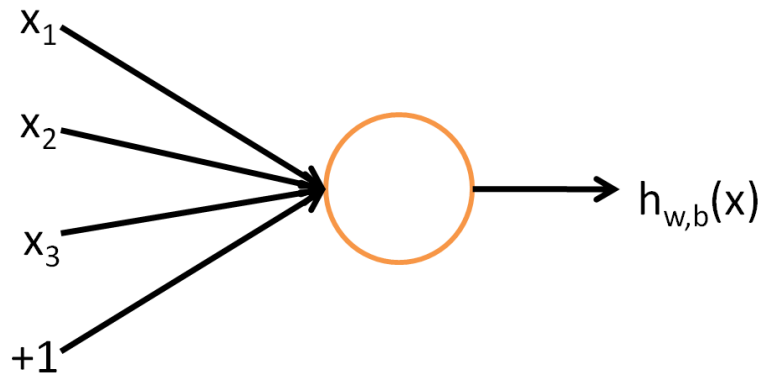


Figure 13 A conceptual overview of a single neuron [23].

Figure 11 shows a conceptual image of a "neuron" and its function as a distinct computational element, accepting inputs x_1, x_2, x_3 , in addition to a +1 intercept component, and subsequently generates a specific output:

$$h_{W,b}(x) = f(W^T x) = f\left(\sum_{i=1}^3 W_i x_i + b\right) \quad (6)$$

where:

$$(f: R \mapsto R) \quad (7)$$

is called the activation function.

These “weights” determine the strength of the connection between the neurons and are adjusted during the training process to minimize the error or cost function. The basic operation of a neural network is forward propagation, which involves passing the input data through the network and computing the output at each layer based on the weights of the connections between the neurons [21].

The output of the final layer is the prediction or decision made by the network. During the training process, the weights of the connections between the neurons are adjusted using an optimization algorithm, such as gradient descent, to minimize the error or cost function [21]. This process is known as backpropagation, where the error is propagated backward through the network to adjust the weights of the connections. Neural networks have been used in a wide range of applications, including image recognition, speech recognition, natural language processing, and many others. They are particularly powerful for tasks that involve complex, non-linear relationships between the input and output, where traditional machine learning models may not perform well.

Each neuron in the MLP receives input from the previous layer, applies a weighted sum of the inputs, and passes the result through an activation function to produce the output. The weights and biases of the neurons are learned during the training process using an optimization algorithm that minimizes a cost function [21].

CONVOLUTION NEURAL NETWORK

A Convolutional Neural Network (CNN) is a type of neural network that is primarily used for image processing and analysis. It is designed to learn spatial hierarchies of features automatically and adaptively from input images [21].

The basic building blocks of a CNN are convolutional layers, pooling layers, and fully connected layers.

Convolutional layers apply filters (also known as kernels or weights) to input images to extract features such as edges, corners, or textures. These filters slide over the input image in a series of matrix multiplications and generate a feature map which highlights the most important features of the input image. Pooling layers are used to reduce the spatial dimensions of the feature maps by taking the maximum or average value of a small region (usually a 2x2 pixel area) of the feature map. Fully connected layers are used to map the learned features to class scores or probabilities. The fully connected layers take the output from the convolutional and pooling layers and apply a traditional feedforward neural network to it [21].

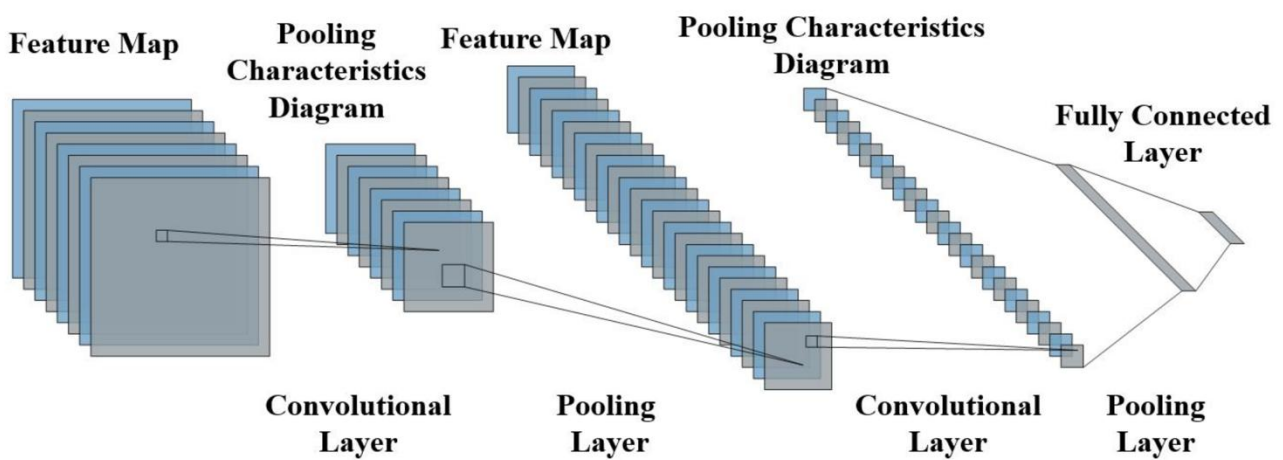
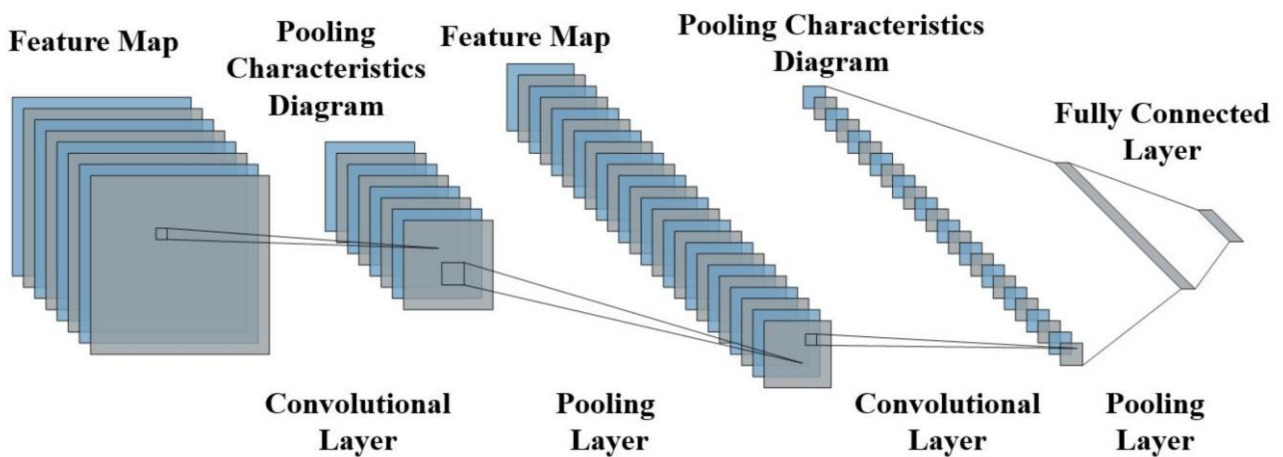


Figure 14 An architectural overview of a convolutional neural network. Taken from [69]



CNNs have shown promising results in image classification, object detection, face recognition, and other computer vision tasks. They can be trained using backpropagation, which adjusts the weights of the network to minimize the difference between the predicted outputs and the ground truth labels.

DATA QUALITY AND MODEL ASSESSMENT

In machine learning, data quality refers to the degree to which data is accurate, complete, relevant, and consistent with the problem being solved. High-quality data is essential for building accurate and reliable machine-learning models.

The key factors that determine data quality are:

- **Accuracy:** The accuracy of the data refers to the degree to which it is free from errors or mistakes. Accurate data is necessary for building reliable machine-learning models.
- **Repeatability:** refers to the ability to reproduce the data consistently, either by collecting the same data again or by using the same process to generate the data. If the data is not repeatable, it can lead to inconsistencies and errors in the machine-learning model. Therefore, it is important to ensure that the data collection and generation process is well-documented and standardized to ensure repeatability.
- **Completeness:** The completeness of the data refers to the degree to which it includes all relevant information. Incomplete data can lead to bias and inaccurate results.
- **Consistency:** The consistency of the data refers to the degree to which it is uniform and consistent across different sources and time periods. Inconsistent data can lead to errors and biases in the machine learning model.
- **Relevance:** The relevance of the data refers to the degree to which it is applicable to the problem being solved. Irrelevant data can lead to inaccurate and biased results.
- **Timeliness:** The timeliness of the data refers to the degree to which it is up-to-date and reflects the current state of the problem being solved. Outdated data can lead to inaccurate and irrelevant results.

Assessing the data quality involves several steps, including:

- **Data Cleaning:** Data cleaning involves identifying and correcting errors in the data, such as missing values, incorrect data types, and outliers.
- **Data Preprocessing:** Data preprocessing involves transforming the data into a format that can be easily understood by the machine learning algorithm. This includes tasks such as scaling, normalization, and feature selection.
- **Exploratory Data Analysis (EDA):** EDA involves visualizing and analyzing the data to identify patterns, trends, and relationships between variables. This can help to identify potential data quality issues, such as inconsistent or missing data.
- **Statistical Analysis:** Statistical analysis involves applying statistical methods to the data to identify patterns and relationships and to test hypotheses. This can help to assess the quality of the data and identify potential sources of bias or errors.
- **Domain Expertise:** Domain expertise involves consulting with experts in the field to ensure that the data is relevant and accurate for the problem being solved. This can help to identify potential limitations or biases in the data.

Assessing data quality is a crucial step in machine learning, as it helps to ensure that the machine learning model is accurate and reliable. It involves a combination of data cleaning, preprocessing, exploratory data analysis, statistical analysis, and domain expertise.

Evaluating a machine learning model involves measuring its performance on a dataset that was not used to train the model. This is done to assess how well the model generalizes to new data and to identify potential problems, such as overfitting or underfitting.

There are several metrics that can be used to evaluate a machine learning model, depending on the problem being solved and the type of model being used. Here are some commonly used evaluation metrics:

- **Accuracy:** Accuracy measures the proportion of correctly classified instances among all instances in the test set.
- **Precision and Recall:** Precision and recall are metrics used in binary classification problems to measure the ability of the model to correctly identify positive instances (precision) and to identify all positive instances (recall).

Accuracy, precision, and recall (sensitivity) are calculated as

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (8)$$

$$Precision = \frac{TP}{TP+FP} \quad (9)$$

$$Recall(sensitivity) = \frac{TP}{TP + FN} \quad (10)$$

where TP is the number of true positives, i.e., data samples/passings that belong to this type and are correctly classified; FP is the number of false positives, i.e., data samples that do not belong to this type and are incorrectly classified as this type; FN is the number of false negatives, i.e., data samples that belong to this type and are incorrectly classified as other types; and TN is the number of true negatives, i.e., data samples that do not belong to this type and are correctly classified as other types.

- **F1-Score:** The F1-score is the harmonic mean of precision and recall and is often used to balance the tradeoff between precision and recall.
- **ROC-AUC:** The Receiver Operating Characteristic - Area Under the Curve (ROC-AUC) measures the tradeoff between true positive rate and false positive rate for different classification thresholds.
- **Mean Squared Error:** Mean Squared Error (MSE) is a commonly used metric to evaluate regression models, which measures the average squared difference between the predicted and actual values.
- **Mean Absolute Error:** Mean Absolute Error (MAE) is another commonly used metric for evaluating regression models, which measures the average absolute difference between the predicted and actual values.

To evaluate a machine learning model, the dataset is typically split into two sets: a training set and a test set. The model is trained on the training set and then evaluated on the test set using one or more of the evaluation metrics described above. It is important to ensure that the test set is representative of the data that the model will encounter in the real world and to avoid any data leakage between the training and test sets. Cross-validation can also be used to evaluate the model on multiple subsets of the data and to obtain a more robust estimate of its performance.

SATELLITE DATA QUALITY CONCEPTS AND REQUIREMENTS

Satellite data quality is an essential aspect of remote sensing applications, as it directly impacts the reliability and usefulness of the information derived from satellite imagery. To ensure the accuracy and precision of satellite data, several concepts and requirements must be considered [24]:

1. **Spatial resolution:** The spatial resolution refers to the smallest ground area that can be distinguished by a satellite sensor. Higher spatial resolution enables the detection of finer details on the Earth's surface, which is crucial for applications such as urban planning, environmental monitoring, and disaster management.
2. **Temporal resolution:** Temporal resolution is the frequency at which a satellite sensor revisits the same location on Earth. The higher temporal resolution allows for more frequent observations, which is particularly important for monitoring rapidly changing phenomena such as wildfires, floods, and crop growth.
3. **Spectral resolution:** Spectral resolution refers to the number and width of the wavelength bands that a satellite sensor can detect. Different spectral bands provide information about various surface properties, such as vegetation health, water quality, and mineral composition. Higher spectral resolution enables more accurate characterization of the Earth's surface features.
4. **Calibration and validation:** Satellite sensors must be regularly calibrated to ensure that their measurements remain accurate and consistent over time. Calibration involves comparing satellite data with ground-based measurements or other reference data sources. Validation is the process of assessing the accuracy of satellite-derived products by comparing them with independent ground-based observations.
5. **Data processing and correction:** Satellite data must be processed and corrected for various sources of error and distortion, such as atmospheric effects, sensor noise, and geometric distortions. Advanced processing algorithms and techniques are employed to minimize these errors and improve the overall data quality.
6. **Data interoperability and standardization:** Satellite data from different sources and sensors should be interoperable and standardized to facilitate comparison and integration for various applications. Standardization ensures that data from different satellites can be easily combined, compared, and used together in a consistent manner.
7. **Level 1 (L1) and Level 2 (L2) Data:** Another important aspect to consider is the distinction between different levels of data processing, particularly L1 and L2 data. L1 data is raw satellite data that has been minimally processed, primarily to correct for sensor-related errors and distortions. In contrast, L2 data has undergone additional processing steps to correct for environmental effects and enhance the data's usability. Understanding the differences between L1 and L2 data can help users better assess the quality and suitability of satellite data for their specific applications.

By adhering to these concepts and requirements, satellite data quality can be optimized, ensuring that the information obtained from remote sensing is reliable and valuable for a wide range of tasks and applications.

SATELLITE IMAGERY: RESOLUTION VS. ACCURACY

Resolution and accuracy are two important properties of satellite data in remote sensing applications, and they can have a significant impact on the quality and usefulness of the data.

Resolution refers to the level of detail or spatial granularity in the data, typically measured in terms of the size of the smallest feature that can be resolved by the sensor. Higher-resolution data can provide more detailed information about the features on the Earth's surface but may also require more processing and storage resources. Lower-resolution data, on the other hand, may have less detail but can cover larger areas and be more easily processed.

Accuracy, on the other hand, refers to the degree of correspondence between the measured data and the true value or state of the feature being measured. High-accuracy data is more reliable and can provide more precise information about the features on the Earth's surface but may require more calibration and validation efforts. Lower-accuracy data, on the other hand, may have more uncertainty and errors and may require more sophisticated data processing and analysis techniques.

The impact of resolution vs. accuracy depends on the specific application and the types of features being measured. For example, in agriculture applications, high-resolution data can provide more detailed information about crop growth and health, while high-accuracy data can provide more precise information about soil moisture and nutrient levels. In urban planning applications, high-resolution data can provide more detailed information about building and infrastructure location and height, while high-accuracy data can provide more precise information about surface elevation and terrain characteristics.

In general, it is important to balance the trade-off between resolution and accuracy depending on the specific needs of the application. In some cases, it may be necessary to compromise on one or the other to achieve the best overall results. For example, in some applications, it may be more important to have accurate data over a large area, even if the resolution is relatively low. In other applications, high-resolution data may be more important, even if the accuracy is lower [25] [26].

SPECTRAL RANGE

Spectral range refers to the range of wavelengths or frequencies of electromagnetic radiation that a satellite sensor can detect and measure. Satellite sensors typically operate in one or more spectral ranges, depending on the type of sensor and the application [27].

In remote sensing applications, the spectral range of a satellite sensor is important because different types of land cover, vegetation, water bodies, and other features have unique spectral signatures or reflectance patterns in different parts of the electromagnetic spectrum. By analyzing the spectral characteristics of the data collected by the sensor, it is possible to identify and map different types of features on the Earth's surface [27].

The spectral range of a satellite sensor is typically divided into several bands, each with a specific range of wavelengths or frequencies. For example, the Landsat series of satellite sensors operate in several spectral bands, including visible, near-infrared, shortwave infrared, and thermal infrared bands. Other satellite sensors, such as the MODIS and VIIRS sensors, operate in a broader range of spectral bands, including visible, near-infrared, shortwave infrared, and thermal infrared, as well as other specialized bands, such as atmospheric sounding bands.

The choice of spectral range for a satellite sensor depends on the specific application and the types of features that need to be detected and measured. For example, a sensor designed for vegetation monitoring might have a high-resolution near-infrared band, which is sensitive to chlorophyll absorption and can provide, e.g., information on vegetation health and productivity [27].

SIGNAL-TO-NOISE RATIO (SNR)

Signal-to-Noise Ratio (SNR) is a measure of the quality of a signal relative to the amount of noise that is present in the signal. In the context of satellite data, SNR refers to the ratio of the strength of the signal (i.e., the information of interest, such as the data collected by a satellite sensor) to the level of background noise (i.e., any unwanted signals that interfere with the data) [28].

In satellite data, SNR can be affected by various factors, such as atmospheric conditions, sensor noise, and interference from other sources. A high SNR means that the signal is stronger than the noise, and the data is more reliable and accurate. On the other hand, a low SNR means that the noise is stronger than the signal, and the data may be corrupted or distorted.

To improve the SNR of satellite data, various techniques can be used, such as signal filtering, noise reduction, and data fusion. For example, in remote sensing applications, the SNR of satellite imagery can be improved by applying atmospheric correction algorithms, which remove the effects of atmospheric scattering and absorption on the signal. In communication applications, the SNR of satellite signals can be improved by using error-correcting codes, which can detect, and correct errors caused by noise.

APPLICATION USER INTERFACE

Application user interfaces (APIs) serve as a vital component in modern software, enabling different systems to communicate and exchange information efficiently. Developers can access features, functions, or data provided by other systems through APIs without delving into underlying implementation details. For remote sensing applications, such as methane emission estimation, APIs play a crucial role in accessing and processing large-scale geospatial datasets from satellite data providers.

APIs offer numerous benefits, such as ease of integration, faster development cycles, and improved scalability, allowing users to access data from multiple sources, perform complex calculations, and present results in a unified, user-friendly format. Two primary user interfaces and APIs for accessing and processing Sentinel data are the Google Earth Engine API and the Sentinel API.

GOOGLE EARTH ENGINE API

The Google Earth Engine API is a cloud-based platform that enables users to access and process large-scale geospatial datasets, such as those from Sentinel-2 and Sentinel-5P satellites. In this thesis, the Google Earth Engine API is employed to access the satellite data efficiently and perform various preprocessing and enhancement steps, such as cloud masking, atmospheric correction, and application of the Normalized Burn Ratio (NBR) index.

The main advantages of using the Google Earth Engine API include its ability to process large volumes of data quickly and efficiently, as well as its support for various programming languages, such as Python and JavaScript. Additionally, the API offers access to a vast repository of satellite imagery and tools for spatial and temporal analysis. This makes it a powerful and versatile platform for researchers working on remote sensing applications, such as methane emission estimation from flaring activities in the oil and gas industry.

However, there are some limitations associated with using the Google Earth Engine API. These may include restrictions on data access and usage for non-academic or commercial purposes, as well as potential challenges in learning and mastering the API for users without prior programming experience.

The documentation for the Google Earth Engine API can be found at [29].

COPERNICUS OPEN ACCESS HUB

The Copernicus Open Access Hub API, or API Hub, is a specialized service for users who want to automate data retrieval and other tasks using a scripting interface. API Hub shares user credentials with the Scientific Data Hub (SciHub), granting new users API Hub access a week after registering on SciHub. Moreover, changes made to a user's SciHub account, such as password, email, or country, are synchronized with the API Hub account within a week.

For those interested in creating their own scripts to interact with the API, the User Guide provides detailed instructions and guidance. By leveraging the capabilities of the Copernicus Open Access Hub API, users can automate and streamline their data retrieval workflows, enhancing the efficiency and scalability of their satellite data analysis tasks.

SUMMARY

The overall ambition is to enable a user easy access to software that uses satellite data to measure the hourly emission of methane from a emission source. A conceptual idea of such solution can be seen in the following Figure where the user finds and marks the area of interest then analyses it and outputs the estimated “emission mass flow” from the located methane source in ideally real-time.

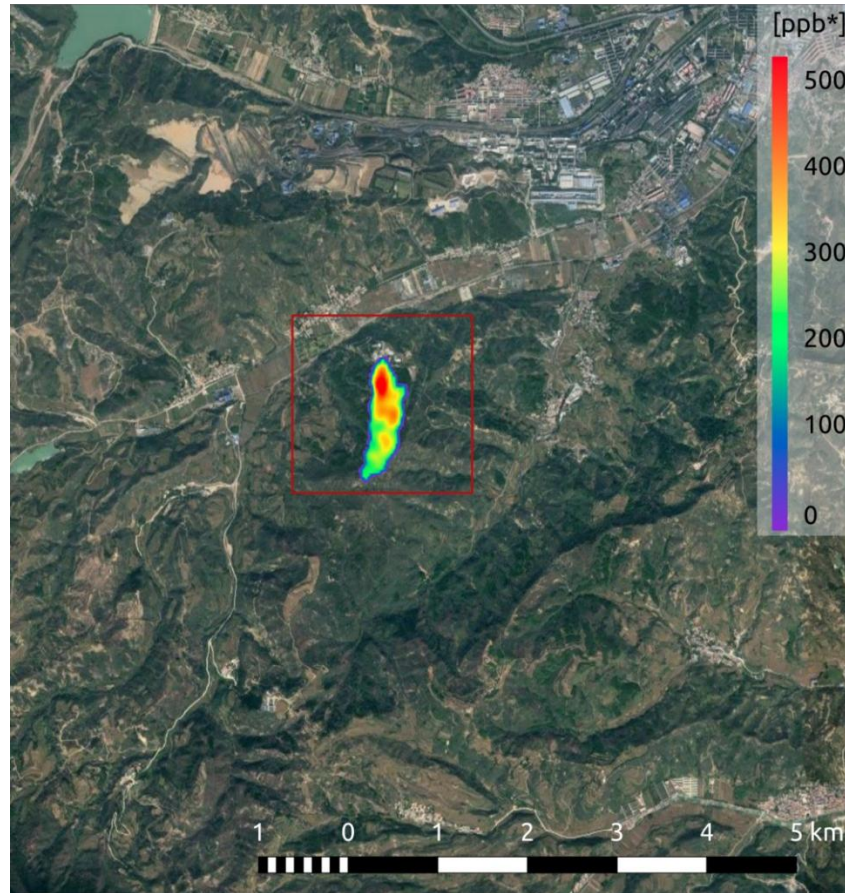


Figure 15 Conceptual design of a methane emission tool

To enable this concept, there are several types of Machine learning approaches. Convolutional Neural Networks (CNNs) can be used to process satellite data in several ways, depending on the specific application and the type of satellite data being used. Here are some examples:

- **Image Classification:** CNNs can be trained to classify satellite images into different categories, such as land use types (e.g., forests, urban areas, water bodies), crop types, or weather patterns. The CNN learns to extract relevant features from the satellite images, such as texture, color, and shape, and use them to classify the images. This can be useful for applications such as environmental monitoring, agriculture, and weather forecasting.
- **Object Detection:** CNNs can be used to detect and locate specific objects within satellite images, such as buildings, vehicles, or ships. This is typically done using a technique called object detection, which involves training a CNN to predict the coordinates of bounding boxes around the objects of interest. This can be useful for applications such as urban planning, disaster response, and maritime surveillance.
- **Change Detection:** CNNs can be used to detect changes in satellite images over time, such as changes in land cover or changes in infrastructure. This is typically done by comparing pairs of satellite images taken at different times and training a CNN to identify the areas where significant changes have occurred. This can be useful for applications such as environmental monitoring, urban planning, and infrastructure management.
- **Image Enhancement:** CNNs can be used to enhance the quality of satellite images, such as by removing noise, sharpening edges, or increasing the resolution. This is typically done by training a CNN to learn a mapping

function between low-quality images and high-quality images. This can be useful for applications such as remote sensing, surveillance, and reconnaissance.

Therefore, CNN seems to be a valid starting point for further work on improving the data quality from satellite platforms.

To evaluate the Data Quality and technical suitability of satellite data, the following metrics have been contextualized in a “performance matrix.”

- **Resolution**
 - **Spectral resolution:** Spectral resolution refers to the number and width of spectral bands captured by the satellite sensor. A higher spectral resolution means that more detailed information about the spectral characteristics of the objects can be captured. The spectral resolution is usually measured in nanometers or micrometers.
 - **Radiometric resolution:** Radiometric resolution refers to the sensitivity of the satellite sensor to differences in the intensity of radiation. A higher radiometric resolution means that smaller differences in intensity can be detected.
 - **Spatial resolution:** Spatial resolution refers to the size of the smallest object that can be resolved by the sensor on the satellite. A higher spatial resolution means that smaller objects can be detected and distinguished. Spatial resolution is usually measured in meters and can be an important factor
- **Revisit time:** Revisit time refers to the time it takes for the satellite to pass over the same location on Earth again. A shorter revisit time means that data can be collected more frequently, which can be important for providing mass flow estimation.
- **Geolocation accuracy:** Geolocation accuracy refers to the accuracy with which the satellite can determine the location of a particular object on Earth. A higher geolocation accuracy means that the location of objects can be determined more precisely, which can be important for applications such as mapping an emission source.
- **Data availability/Cost:** Cost refers to the financial cost of accessing the satellite platform.

The following performance matrix has been put forward, ranging from not suitable (red), partly suitable (yellow), and suitable (green) based on the following metrics.

Table 1 Performance Matrix for Satellite Platforms with examples of satellites as input.

Satellite data	Suitable	Partly suitable	Not suitable
Data Availability	Sentinel-2	Landsat-8	GHGSat-MethaneSat
Revisit time	Sentinel-5P	Sentinel-2	Landsat-8
Resolution	Sentinel-2, MethaneSat	GHGSat- Landsat-8	Sentinel-5P
Geolocation accuracy:	Landsat-8	Sentinel-2	Sentinel-5P

This will be further used when evaluating different satellite platforms, and populated in the technical review.

TECHNICAL REVIEW

To establish a solid foundation for the development of an optimal method for estimating emissions from flaring activities, it is essential to examine the current state of the art in the field. This section of the project will conduct a comprehensive review focused on addressing the key questions outlined below:

1. Current methods of emission monitoring
 - a. How is the emission data via satellite being gathered today?
 - b. How is flaring being reported/monitored?
2. What available satellite “state-of-the-art” platforms are operational or/and are available?
3. What key research has been performed on the topic of using Machine Learning to improve the data quality from satellites?

FLARING MONITORING – HOW IS EMISSION FROM FLARING CURRENTLY BEING MONITORED?

Ultrasonic meters are not typically used to directly measure emissions from flaring, but they can be used to monitor the flow rate of gas through a flare stack, which can be used to estimate the amount of gas being flared [30].

Ultrasonic flow meters use sound waves to measure the velocity of gas flowing through a pipe or stack. They work by transmitting sound waves through the gas stream and measuring the time it takes for the waves to travel upstream and downstream [30]. The difference in travel time can be used to calculate the velocity of the gas, which can then be used to calculate the flow rate, as seen in the following illustration.

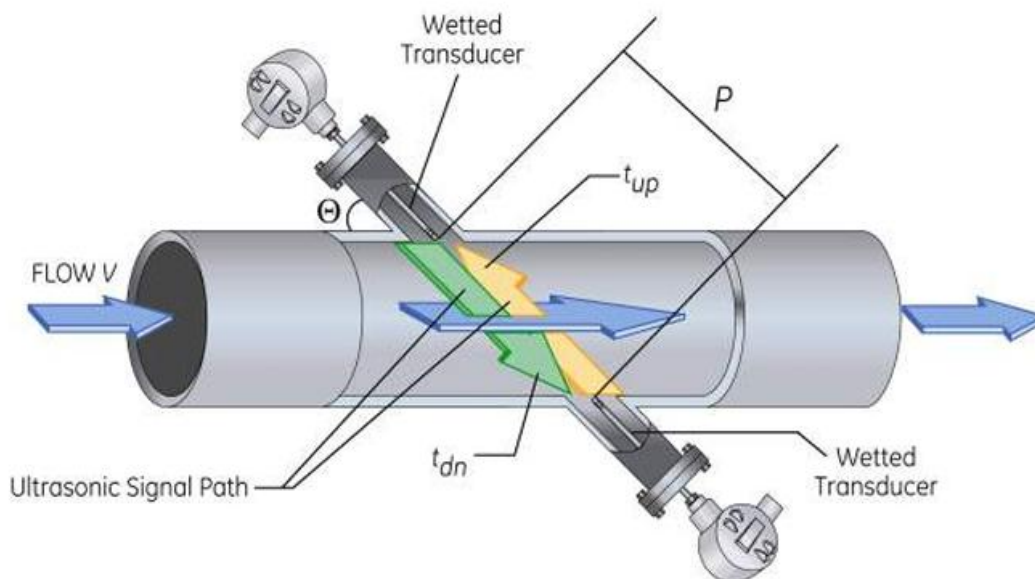


Figure 16 The operating principle of a transit-time-based ultrasonic flowmeter [30].

By measuring the flow rate of gas through a flare stack, operators can estimate the amount of gas being flared and calculate corresponding emissions of greenhouse gases which are key to complying with the current reporting scheme. However, this method is not as accurate as direct measurements of gas composition and concentration, and it may not capture all sources of emissions from the flaring process [30].

It's important to note that ultrasonic meters are just one tool in a suite of technologies used to monitor flaring emissions. In most cases, a combination of on-site measurements and remote sensing technologies is used to provide a comprehensive picture of flaring emissions and ensure compliance with regulatory requirements, but ultrasonic technology is stated as the “best practice.”

EMISSION MONITORING VIA SATELLITES

As stated in the introduction, depending on the respective payload, satellites can be used to monitor emissions by measuring various atmospheric parameters, such as the concentration of pollutants, the temperature of the air, and the concentration of gases in the atmosphere.

One way that satellites can do this is through remote sensing. Satellites equipped with sensors can detect the wavelengths of light that are emitted or reflected by the Earth's surface and atmosphere. By analyzing the patterns of these wavelengths, scientists can infer the concentration of different gases and pollutants in the atmosphere, as seen in the following Figure

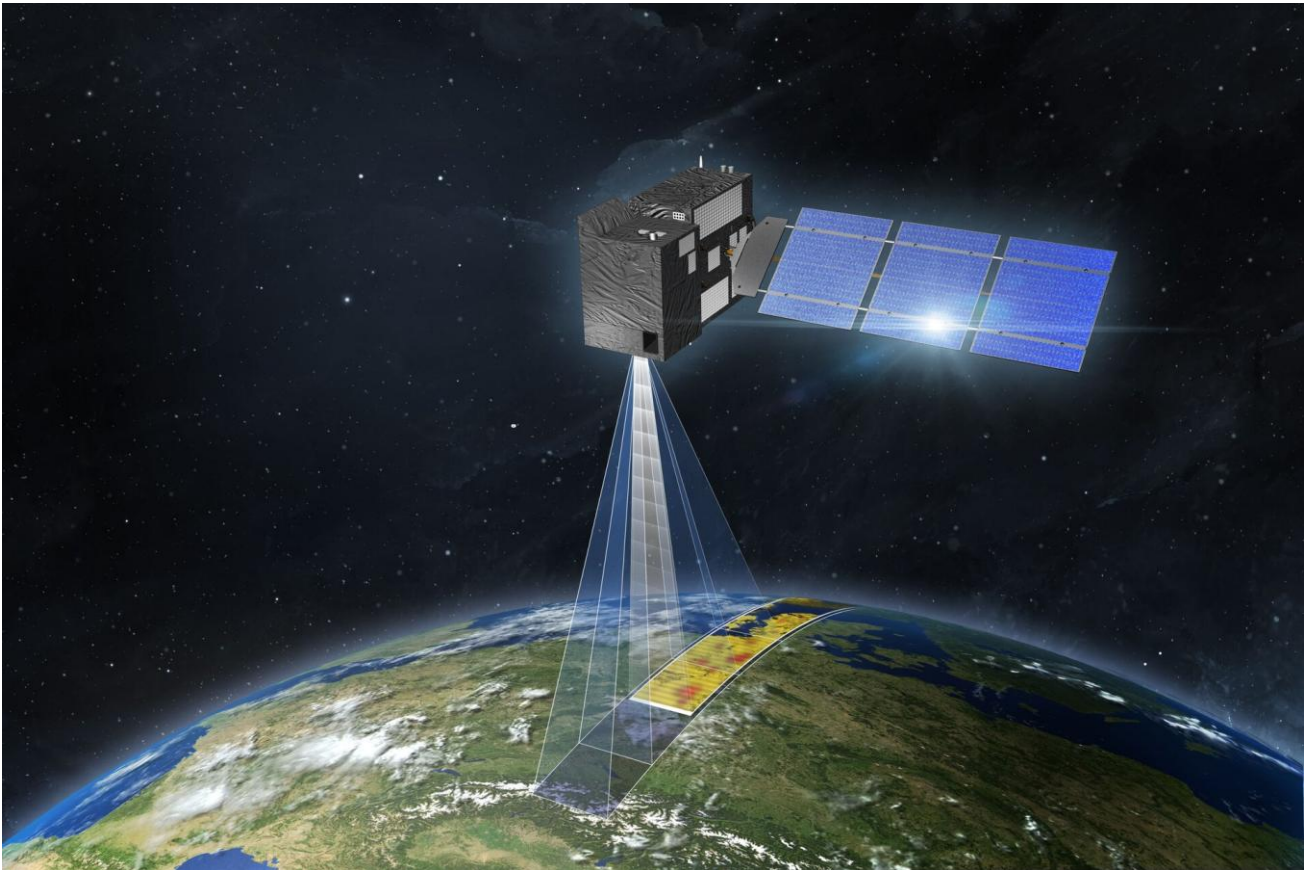


Figure 17 Example of a satellite application is CO₂ measurement from orbit, where the satellite is scanning the earth's surface via its sensory payload [31].

For example, satellites can detect the concentration of greenhouse gases like carbon dioxide and methane by measuring their absorption of certain wavelengths of infrared light. They can also detect the concentration of pollutants like nitrogen dioxide and sulfur dioxide by measuring the wavelengths of light that they absorb or emit [31]. Based on this, it enables a method of measuring the gas concentration of a target gas, usually in part-per-million or parts-er-billion.

Satellites can also use GPS technology to track the movement of emissions sources, such as factories, power plants, and transportation networks. By analyzing changes in GPS signals, scientists can estimate the amount of emissions being released from these sources and track their movements over time.

TECHNICAL REVIEW OF AVAILABLE SATELLITE PLATFORMS

Several companies have already developed satellite-based methane detection technologies, such as TROPOMI, GOSAT, and GHGSat. These technologies use different types of sensors, including hyperspectral sensors, short-wave infrared sensors, and thermal infrared sensors. By combining these sensors, they can detect and quantify methane emissions with high accuracy and spatial resolution. This section is a review of the different satellite platforms and their respective data quality, in addition to an assessment of the availability of satellite data.

SENTINEL-5 AND 2

Sentinel-5 and Sentinel-2 are two distinct Earth observation satellites within the European Space Agency's (ESA) Copernicus Programme. The primary goal of the Copernicus Programme is to monitor Earth's environment and provide valuable information for various applications, such as climate change monitoring, natural resource management, and disaster response.

SENTINEL-5: ATMOSPHERIC MONITORING AND TRACE GAS RETRIEVAL

Sentinel-5 is a satellite mission dedicated to monitoring Earth's atmosphere, primarily focusing on the retrieval of trace gases and atmospheric pollutants. Launched in 2021, Sentinel-5 carries the TROPospheric Monitoring Instrument (TROPOMI), which is a spectrometer capable of measuring sunlight scattered or absorbed by Earth's atmosphere. This instrument covers the ultraviolet, visible, near-infrared, and shortwave infrared spectral ranges [32].

Sentinel-5's primary goal is to provide accurate, high-resolution data on atmospheric constituents, vital for tracking air quality, climate change, and the ozone layer. With TROPOMI's high-resolution data (7x7 km) and global coverage, Sentinel-5 is invaluable for scientists, policymakers, and industries involved in atmospheric monitoring and climate change mitigation [32].

The TROPOMI instrument onboard the Copernicus Sentinel-5 Precursor is a nadir-viewing imaging spectrometer covering wavelength bands between the ultraviolet and the shortwave infrared. The instrument uses passive remote sensing techniques to attain its objective by measuring, at the Top Of Atmosphere (TOA), the solar radiation reflected by and radiated from Earth. The instrument operates in a push-broom configuration (non-scanning), with a swath width of ~2600 km on the Earth's surface. The typical pixel size (near nadir) will be 7x3.5 km² for all spectral bands, except for the UV1 band (7x28 km²) and SWIR bands (7x7 km²) [32].

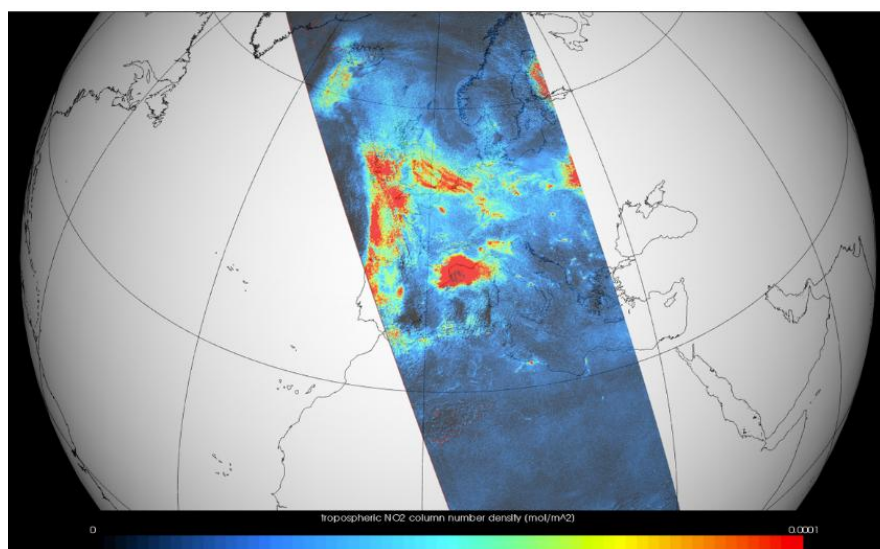


Figure 18 Example of Sentinel-5P NO₂ data visualized [33]

Sentinel-5's high-resolution data and daily global coverage make it a valuable resource for monitoring atmospheric constituents and tracking air quality, climate change, and the ozone layer. This publicly available and free-to-use data benefits scientists, policymakers, and industries involved in environmental monitoring and climate change mitigation.

SENTINEL-2: HIGH-RESOLUTION MULTISPECTRAL IMAGING FOR LAND AND COASTAL MONITORING

The Sentinel-2 mission consists of two satellites, Sentinel-2A (launched in 2015) and Sentinel-2B (launched in 2017). These satellites are designed to provide high-resolution multispectral imagery for land and coastal monitoring, serving a wide range of applications such as agriculture, forestry, urban planning, and disaster management [34].

Equipped with the Multispectral Instrument (MSI), the Sentinel-2 satellites capture data in 13 spectral bands. These bands cover the visible, near-infrared (NIR), and shortwave infrared (SWIR) regions of the electromagnetic spectrum, delivering essential information for various remote sensing tasks. The spatial resolution of Sentinel-2 imagery varies from 10 meters to 60 meters, depending on the specific spectral band. Together, the two satellites achieve a 5-day revisit time at the equator [34].

The high-resolution multispectral data from Sentinel-2 has significantly improved the accuracy and timeliness of land cover and vegetation monitoring. It enables precise detection of changes in land use, assessment of crop health, and identification of natural disasters. Moreover, the imagery is valuable for monitoring coastal and inland waters, offering insights into water quality, sediment transport, and aquatic ecosystems. Figure 19 provides an example of the mission's capabilities.



Figure 19 Acquired on 27 June 2015 at 10:25 UTC (12:25 CEST), just four days after launch, this close-up of France's southern coast from Nice airport (lower left) to Menton (upper right) is a subset from the first image from the Sentinel-2A satellite. This false colour image was processed including the instrument's high-resolution infrared spectral channel [35]

Sentinel-2's data is freely available and provides timely and accurate high-resolution multispectral imagery, benefiting land management, agriculture, disaster response, and environmental monitoring worldwide.

LANDSAT-8

Landsat 8 is a satellite mission operated by NASA and the United States Geological Survey (USGS) that provides global coverage of the Earth's land surface with moderate-resolution multispectral imagery. The satellite was launched on February 11, 2013, and has been in operation since then, collecting data on a continuous basis.



Figure 20 Landsat 8 Satellite Sensor with a 15m resolution [36]

The Landsat 8 satellite, equipped with the Operational Land Imager (OLI) and the Thermal Infrared Sensor (TIRS), captures high-resolution data in multiple spectral bands. The OLI sensor covers visible, near-infrared, shortwave infrared, and panchromatic bands, while the TIRS sensor focuses on thermal infrared bands for temperature measurements. Landsat 8 surpasses its predecessors in terms of improved spectral and radiometric capabilities, offering enhanced signal-to-noise ratio and radiometric calibration. The OLI sensor introduces two additional spectral bands for improved atmospheric correction and cloud detection. With a global coverage cycle of 16 days and a swath width of 185 kilometers, Landsat 8 provides valuable data for various applications, such as land cover mapping, vegetation monitoring, natural resource management, urban planning, and disaster response. This freely available data from the USGS EarthExplorer website enables researchers and analysts to make informed decisions regarding resource management and environmental conservation [37].

The Landsat-8 specification can be viewed in Table 2 below.

Table 2 The landsat-8 specification. Data from [37]

Processing:	Level 1 T- Terrain Corrected	
Spatial Resolution		
Band # and Type	Bandwidth (µm)	Resolution (m)
Band 1 Coastal	0.43 - 0.45	30
Band 2 Blue	0.45 - 0.51	30
Band 3 Green	0.53 - 0.59	30
Band 4 Red	0.63 - 0.67	30
Band 5 NIR	0.85 - 0.88	30
Band 6 SWIR 1	1.57 - 1.65	30
Band 7 SWIR 2	2.11 - 2.29	30
Band 8 Pan	0.50 - 0.68	15
Band 9 Cirrus	1.36 - 1.38	30
Band 10 TIRS 1	10.6 - 11.19	30 (100)
Band 11 TIRS 2	11.5 - 12.51	30 (100)
Pixel Size:	<ul style="list-style-type: none"> • OLI multispectral bands 1-7,9: 30-meters • OLI panchromatic band 8: 15-meters 	

	<ul style="list-style-type: none">• TIRS bands 10-11: collected at 100 meters but resampled to 30 meters to match OLI multispectral bands
Data Characteristics:	<ul style="list-style-type: none">• GeoTIFF data format• Cubic Convolution (CC) resampling• North Up (MAP) orientation• Universal Transverse Mercator (UTM) map projection (Polar Stereographic for Antarctica)• World Geodetic System (WGS) 84 datum• 12 meter circular error, 90% confidence global accuracy for OLI• 41 meter circular error, 90% confidence global accuracy for TIRS 16-bit pixel values

WORLDVIEW-3

WorldView-3 is a commercial satellite mission jointly operated by DigitalGlobe and Maxar, prominent providers of high-resolution satellite imagery and geospatial solutions. Launched on August 13, 2014, the satellite has been actively capturing high-resolution imagery of the Earth's surface for various applications.



Figure 21 Conceptual rendering of Maxar's WorldView-3 satellite [38]

Equipped with a multi-spectral sensor, WorldView-3 can acquire data in 16 spectral bands, including four new bands not present in previous DigitalGlobe satellites. The sensor offers a ground resolution of 31 centimeters for panchromatic imagery and 1.24 meters for multispectral imagery, enabling detailed mapping and analysis of surface features. One notable feature of WorldView-3 is its capability to collect data in the shortwave infrared (SWIR) part of the electromagnetic spectrum. This allows the satellite to penetrate atmospheric conditions like haze, smoke, and dust, enabling the detection and identification of materials based on their distinctive spectral signatures. The SWIR bands are particularly valuable for applications such as mineral exploration, environmental monitoring, and military uses. WorldView-3 also boasts high geolocation accuracy, with a pointing accuracy of under 3 meters and a knowledge accuracy of under 1 meter [39]. This precision facilitates precise mapping and geospatial analysis of surface features and structures. The complete specification for WorldView-3 can be seen in the following table.

Table 3 WorldView-3 sensor specification. Data from [39]

Design	Specifications
Sensor resolution	Panchromatic 0.3 m 1.24m multispectral Short-wave 3.7m CAVIS 30 m
Sensor bands	
Panchromatic	450 – 800 nm
Multispectral	8 bands
Costal	400-500 nm
Blue	450-510 nm
Green	510 -580 nm
Yellow	585-625 nm
Red	630 – 690 nm
Near-IR	770 – 850 nm
Near-IR2	860 – 1040 nm

The satellite has a revisit time of 1-3 days, depending on the location and imaging requirements. The data is available for commercial use through DigitalGlobe's cloud-based platform and can be used for a wide range of applications, including urban planning, natural resource management, defense and intelligence, and disaster response.

Maxar's methane mapping algorithm can automatically identify and quantify methane emissions of greater than 1,000 kg per hour with a wind speed of 10 miles per hour. Maxar's research notes the visual inspection of the methane mapping layer can identify methane plumes down to much lower emission rates of less than 50 kg per hour with 10 miles per hour wind speed. The actual minimum detectable quantity is influenced by several factors, including:

- Mass flow rate (rate of methane emissions in kg per hour)
- Flow velocity (wind speed for leaks, forced flow for mine vents)
- Point source (e.g., pipeline leak) versus distributed source (e.g., landfill)
- Spectral clutter of the imagery
- Atmospheric stability (turbulence)

WorldView-3's high-resolution satellite imagery, operated by DigitalGlobe and Maxar, offers detailed mapping and analysis capabilities. However, the data is not freely available to the public. Commercial users can access it through DigitalGlobe's platform for various applications including urban planning, natural resource management, defense, and disaster response. Additionally, the European Space Agency (ESA) offers limited access to WorldView-3 data, requiring an application process for eligible users.

GHGSAT- METHANESAT

GHGSat-MethaneSAT is a joint satellite mission between GHGSat, a Canadian company that specializes in satellite-based greenhouse gas monitoring, and EDF, a French energy company. The mission is focused on detecting and monitoring methane emissions from industrial sources worldwide, with a goal of reducing methane emissions by up to 45% by 2025. GHGSat has deployed its microsatellite "Iris" along with six other satellites in orbit. With expertise in high-resolution remote sensing of greenhouse gases, GHGSat aims to reduce methane emissions by up to 45% by 2025. Their methane data products, evaluated under the CSDA program, contribute to the mission's objectives, including the measurement of emissions from individual oil and gas wells and other small point sources [40].

The satellite is equipped with a state-of-the-art hyperspectral sensor that can detect and measure methane concentrations in the atmosphere with high precision and accuracy. The sensor has a resolution of 25 meters, which means it can detect methane emissions from individual facilities and pinpoint the source of the emissions. GHGSat-MethaneSAT uses a hyperspectral sensor that operates in the infrared spectrum to detect and measure methane concentrations in the atmosphere. The sensor is specifically tuned to detect the unique spectral signature of methane, which allows it to distinguish methane from other gases in the atmosphere. The sensor operates in the 1.6 to 1.7-micron range, which is where methane absorbs infrared radiation. This allows GHGSat-MethaneSAT to detect and measure methane emissions with high precision and accuracy.

GHGSat-MethaneSAT is designed to operate in a sun-synchronous orbit at an altitude of 500 kilometers, which allows it to capture high-resolution imagery of the Earth's surface with a revisit time of about once every three days. The satellite is also equipped with advanced data processing and analysis capabilities that enable it to produce accurate and timely methane emissions data for a wide range of industrial sources, including oil and gas operations, landfills, and agriculture [40].

One of the key features of GHGSat-MethaneSAT is its ability to detect methane emissions from smaller sources that are typically overlooked by traditional monitoring methods. This includes emissions from leaks in pipelines and storage tanks, as well as emissions from smaller facilities that are not required to report their emissions under current regulations.

Figure 18 shows the application of the data captured via the GHGSat from an oil and gas facility in which it estimated a mass flow rate of around $190 \frac{kg}{hr}$.

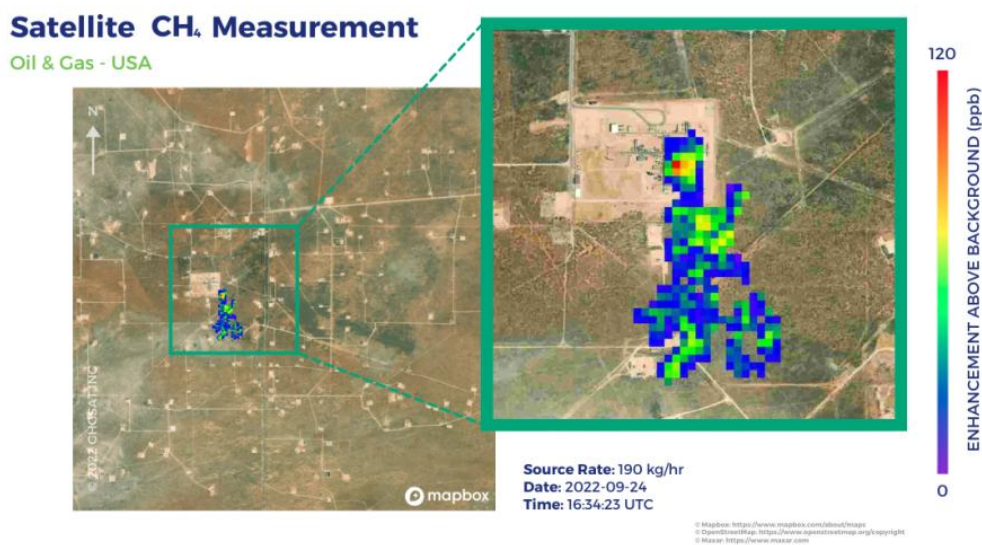


Figure 22 High-resolution satellite measurement by GHGSat of methane emission from an oil and gas facility in New Mexico on September 24, 2022. Taken from [41]

Another case was that a methane plume was spotted near Jordan's capital Amman. GHGSat attributed the plume to the waste sector and estimated its emissions rate at 4,876 kilograms an hour [42].

GHGSat-MethaneSAT provides accurate and timely methane emissions data for various industries. Access to its high-resolution imagery and advanced data processing is not publicly available, and is behind a paywall.

VISIBLE INFRARED IMAGING RADIOMETER SUITE (VIIRS)

The Visible Infrared Imaging Radiometer Suite (VIIRS) is a key instrument aboard the National Oceanic and Atmospheric Administration's (NOAA) Suomi NPP and NOAA-20 satellites. It is designed to provide high-quality global observations of the Earth's surface, atmosphere, and cloud cover and has been operational since 2011 [43].

Technical specifications of VIIRS:

- **Spectral Bands:** VIIRS has 22 spectral bands, including 14 in the reflective solar region (RSB) and 8 in the thermal emissive region (TEB).
- **Spatial Resolution:** VIIRS has five different spatial resolutions, ranging from 375 meters to 7500 meters, depending on the spectral band.
- **Swath Width:** The VIIRS sensor has a swath width of approximately 3000 km, which allows for broad coverage of the Earth's surface in a single pass.
- **Radiometric Resolution:** VIIRS provides a radiometric resolution of 12 bits, which allows for more accurate measurements of radiation and energy.
- **Thermal Sensitivity:** VIIRS has a thermal sensitivity of less than 0.05 K, which enables it to detect very small differences in temperature.
- **Data Products:** VIIRS data is used to generate a variety of products, including sea surface temperature, vegetation index, cloud properties, atmospheric temperature and moisture profiles, and ocean color.

The VIIRS instrument is a state-of-the-art imaging radiometer that provides high-quality data for a wide range of applications, including weather forecasting, environmental monitoring, and climate research. Its high spectral and spatial resolution, along with its advanced radiometric and thermal sensitivity, make it a valuable tool for scientists and researchers studying the Earth's climate and environment.

Table 4 Characteristics of the Nine Visible Infrared Imaging Radiometer Suite (VIIRS) Spectral Bands Collecting Data at Night [44]

Band Designation	Spectral Range	Bandpass (μm)	Band Center (μm)	Lmin Requirement ($\text{W}/(\text{m}^2 \cdot \mu\text{m} \cdot \text{sr})$)	Lmax Requirement ($\text{W}/(\text{m}^2 \cdot \mu\text{m} \cdot \text{sr})$)
DNB	Panchromatic	0.5–0.9	0.7	$3.0\text{E}^{-5} \text{ W}/(\text{m}^2 \cdot \text{sr})$	$200 \text{ W}/(\text{m}^2 \cdot \text{sr})$
M7	Near-infrared	0.843–0.881	0.862	3.4	349
M8	Near-infrared	1.225–1.252	1.2385	3.5	164.9
M10	Short-wave IR	1.571–1.631	1.601	1.1	71.2
M12	Mid-wave IR	3.598–3.791	3.6945	0.0078	2.84
M13	Mid-wave IR	3.987–4.145	4.066	0.00216	406
M14	Long-wave IR	8.407–8.748	8.5775	0.373	19.5
M15	Long-wave IR	10.234–11.248	10.741	0.729	17.1
M16	Long-wave IR	11.405–12.322	11.865	0.876	14.5

VIIRS can be used to monitor methane emissions from flaring in oil and gas production facilities. By detecting and measuring the amount of thermal radiation emitted by the flare, VIIRS can provide information on the amount of methane that is being released. To monitor methane emissions using VIIRS, researchers typically use a technique called thermal infrared remote sensing. This involves measuring the thermal radiation emitted by the flare in the 4.5-5.5 micron spectral range, which is sensitive to methane emissions [45].

VIIRS can detect and measure methane emissions using its thermal emissive bands. These bands are sensitive to the thermal radiation emitted by hot objects, including flares. The VIIRS instrument measures the temperature of the flare and the surrounding area, which allows for the calculation of the flare's radiative power. The radiative power can be used to estimate the amount of methane that is being released during the flaring activity, as seen in the Figure below.

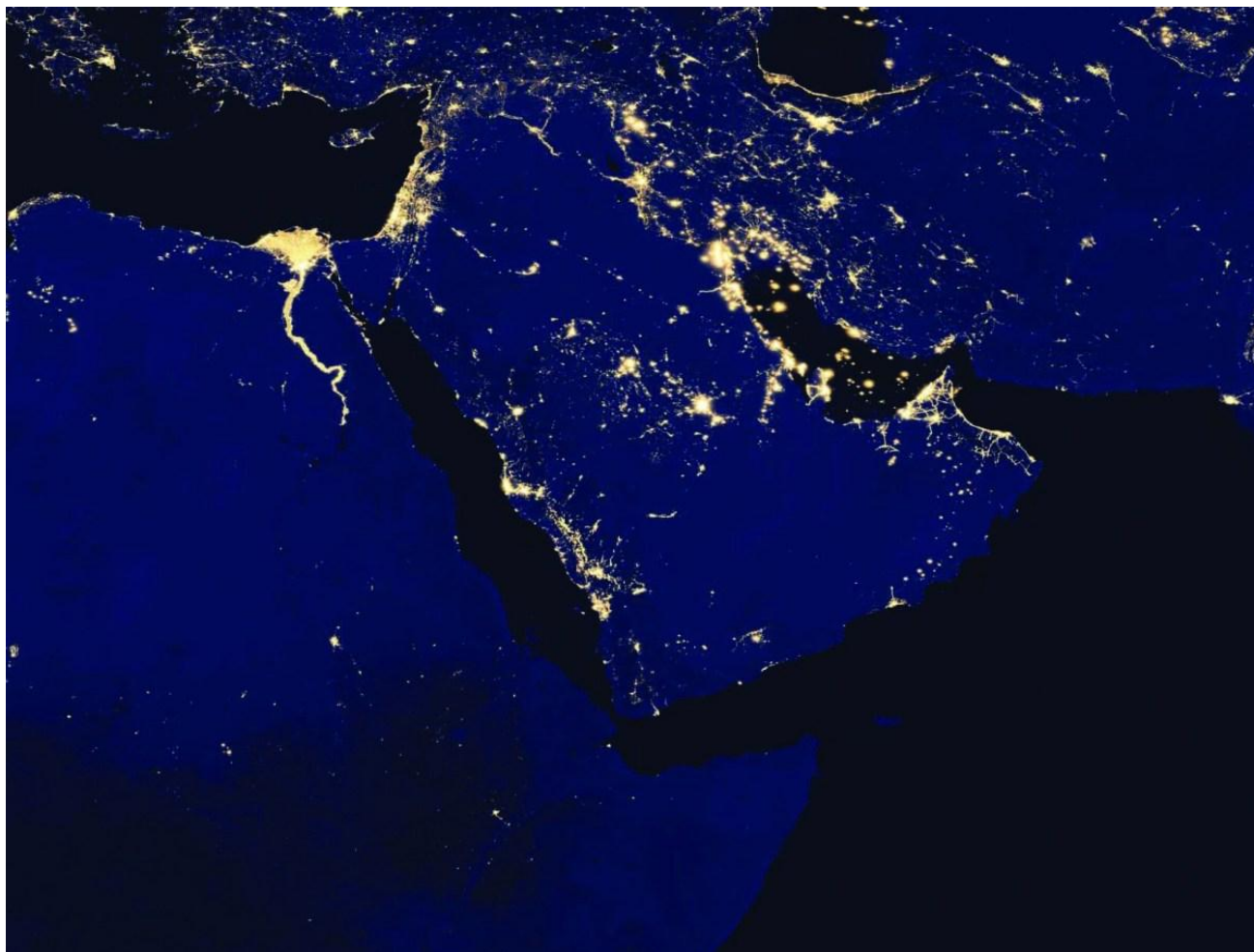


Figure 23 They are using VIIRS to pinpoint areas where fuel production is leading to gas flaring, the practice of burning off natural gas that comes to the surface with crude oil, and to estimate the amount of gas released at those sites [45].

In addition to monitoring methane emissions from flaring, VIIRS can also be used to detect and monitor methane leaks from oil and gas production facilities. Methane leaks can occur during the production, storage, and transportation of oil and gas. VIIRS can detect these leaks by measuring the thermal radiation emitted by the methane gas. This information can be used to identify and locate sources of methane emissions and to develop strategies for reducing emissions.

The data acquired from VIIRS is publicly available, and available at NASA's Earthdata platform.

LITERATURE REVIEW OF RELEVANT RESEARCH

The purpose of this section is to present an overview of the “state-of-the-art” research on methods for improving the data quality from satellites. The emphasis of this literature review is.

1. General improvements in data quality – improved sensors and satellites
2. Improvements based on Machine Learning and existing datasets
3. Specific work on the development of data-driven methods for improving methane emission from satellites

GENERAL IMPROVEMENTS IN DATA QUALITY

This section highlights the work done on improving the data quality from satellite measurements. Independent verification of satellite performance is important for identifying large greenhouse gas point sources because it ensures that the data collected by satellites is accurate and reliable. This is necessary for acceptance and use by policymakers and stakeholders who rely on this information to make informed decisions about mitigation efforts. Existing satellite validation efforts largely focus on the consistency of quantification estimates across satellites and methods or rely on internal and generally unpublished controlled methane release testing. However, independent blind testing has become commonplace for terrestrial and airborne methane sensing technologies, but until now, there has been no such testing of the methane detection and quantification capabilities of satellites due to the expense and logistical challenges of performing releases at scales visible from space.

METHANE RETRIEVED FROM TROPOMI: IMPROVEMENT OF THE DATA PRODUCT AND VALIDATION OF THE FIRST 2 YEARS OF MEASUREMENTS

This work [46] describes the main improvements made to retrieve CH₄ from TROPOMI measurements using the full-physics approach, which is a method of retrieving methane (CH₄) data from TROPOMI measurements that consider the physical properties of the atmosphere and the satellite instrument. This approach uses a radiative transfer model to simulate the interaction between sunlight and atmospheric gases, which allows for more accurate retrieval of CH₄ concentrations. In contrast, a proxy approach assumes that light path modifications due to scattering in the atmosphere are the same for both CH₄ and another gas (in this case, CO₂), which To validate the methane (CH₄) data retrieved from TROPOMI, independent ground-based CH₄ measurements from the Total Carbon Column Observing Network (TCCON) was used as a reference. Thirteen different TCCON stations located in North America, East Asia, Europe, and Oceania were used for the validation. The validation was conducted by comparing the TROPOMI CH₄ data with the TCCON reference dataset being used as a proxy for CH₄.

The findings presented in this work provide valuable insights into the accuracy and reliability of methane (CH₄) data retrieved from TROPOMI measurements. The improvements made to the retrieval process have resulted in a more accurate and precise CH₄ data product, which can be used to better understand the sources and sinks of CH₄ emissions. This is particularly important given that CH₄ is a potent greenhouse gas with a global warming potential of more than 80 times higher than that of carbon dioxide (CO₂). By accurately measuring CH₄ concentrations, researchers can better assess the impact of anthropogenic activities such as agriculture and fossil fuel use on climate change. Additionally, the long-term record provided by TROPOMI measurements allows for ongoing monitoring of CH₄ emissions and their impact on the environment. Overall, these findings contribute to our understanding of the global carbon cycle and help inform climate change mitigation strategies.

SINGLE-BLIND VALIDATION OF SPACE-BASED POINT-SOURCE DETECTION AND QUANTIFICATION OF ONSHORE METHANE EMISSIONS

This work published in Nature in March 2023 outlines a first step toward ongoing, operational blind testing of satellites quantifying methane point sources. The authors conducted a desert-based test where five independent teams analyzed data from one to five satellites each for a total of 11 satellite observations which provides key characteristics of each participating satellite [47]

Table 5 Key characteristics of each participating satellite constellation, from lowest to highest swath width, which is roughly proportional to an instrument’s minimum methane detection limit [47].

Satellite	Coverage	Constellation size	Swath (km)	~ Revisit time (per satellite) (days)	Data availability
GHGSat-C2	Targeted	5 (C1–C5)*	12	14	Commercial
WorldView 3	Targeted	1	13.1	4.5 [†]	Commercial
PRISMA	Targeted	1	3	7	Public
Landsat-8	Global	1	185	16	Public
Sentinel-2	Global	2	290	10	Public

Global coverage refers to a configuration that passively covers most of Earth’s surface over some number of orbits, while targeted coverage refers to a “point-and-shoot” instrument that must be pointed to a particular location.

*GHGSat-C3–C5 were launched after the conclusion of testing.

†For best resolution within 20° off nadir. WorldView 3 has 1-day revisit time at lower guaranteed resolution.

They correctly identified 71% of all emissions and concluded that significantly more blind testing is needed to ensure rapid uptake and trust, as seen in Figure 24.

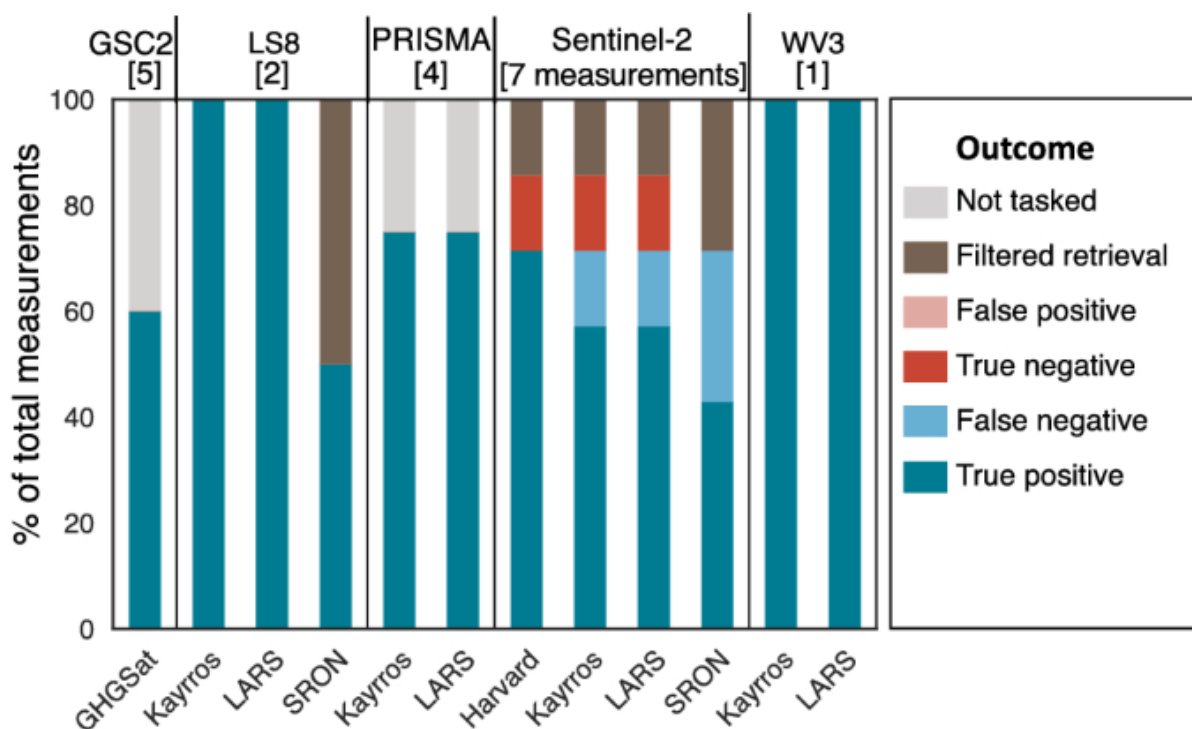


Figure 24 Detection performance by satellite and team. Total number of measurements listed in brackets. For each satellite, most teams correctly detected most emissions as true positives or true negatives (correctly identified non-emissions). In some cases, e.g. two GHGSat-C2 (GSC2) overpasses, the satellite was not tasked and collected no data. In others, e.g. one SRON retrieval of Landsat 8 (LS8), no retrieval was attempted due to image clipping concerns or excessive cloud cover. No teams produced false positives, in which satellites detected methane when none was released [47]

The study provides important insights into the validation of space-based point-source detection and quantification of onshore methane emissions and highlights the need for further research in this area. The findings of this study can be used by policymakers and stakeholders for mitigation efforts in several ways.

1. The study provides evidence that satellites can be used as a tool for identifying large greenhouse gas point sources for mitigation.
2. This work highlights the importance of independent verification of satellite performance to ensure that the data collected by satellites is accurate and reliable.
3. The authors suggest that further blind testing is needed to ensure rapid uptake and trust in satellite-based methane detection and quantification capabilities.

Finally, the study provides a framework for ongoing, operational blind testing of satellites quantifying methane point sources, which can help policymakers and stakeholders make informed decisions about mitigation efforts. Overall, the findings of this study can inform policy-making and stakeholder decision-making related to greenhouse gas emissions reduction efforts which could be directly transferable to flare monitoring.

SATELLITE-BASED SURVEY OF EXTREME METHANE EMISSIONS IN THE PERMIAN BASIN

The basis for this work is that little was known about individual contributors to methane emissions in the Permian basin, which is responsible for almost half of the methane emissions from all U.S. oil- and gas-producing regions. This lack of knowledge made it difficult to develop effective mitigation strategies. Therefore researchers used data from three new hyperspectral satellite missions, each one carrying an imaging spectrometer as payload. These are two versions of the AHSI onboard the GF5 and ZY1 platforms (China, launched in May 2018 and September 2019, respectively) and the core instrument onboard PRISMA (Italy, launched in March 2019 [48]).

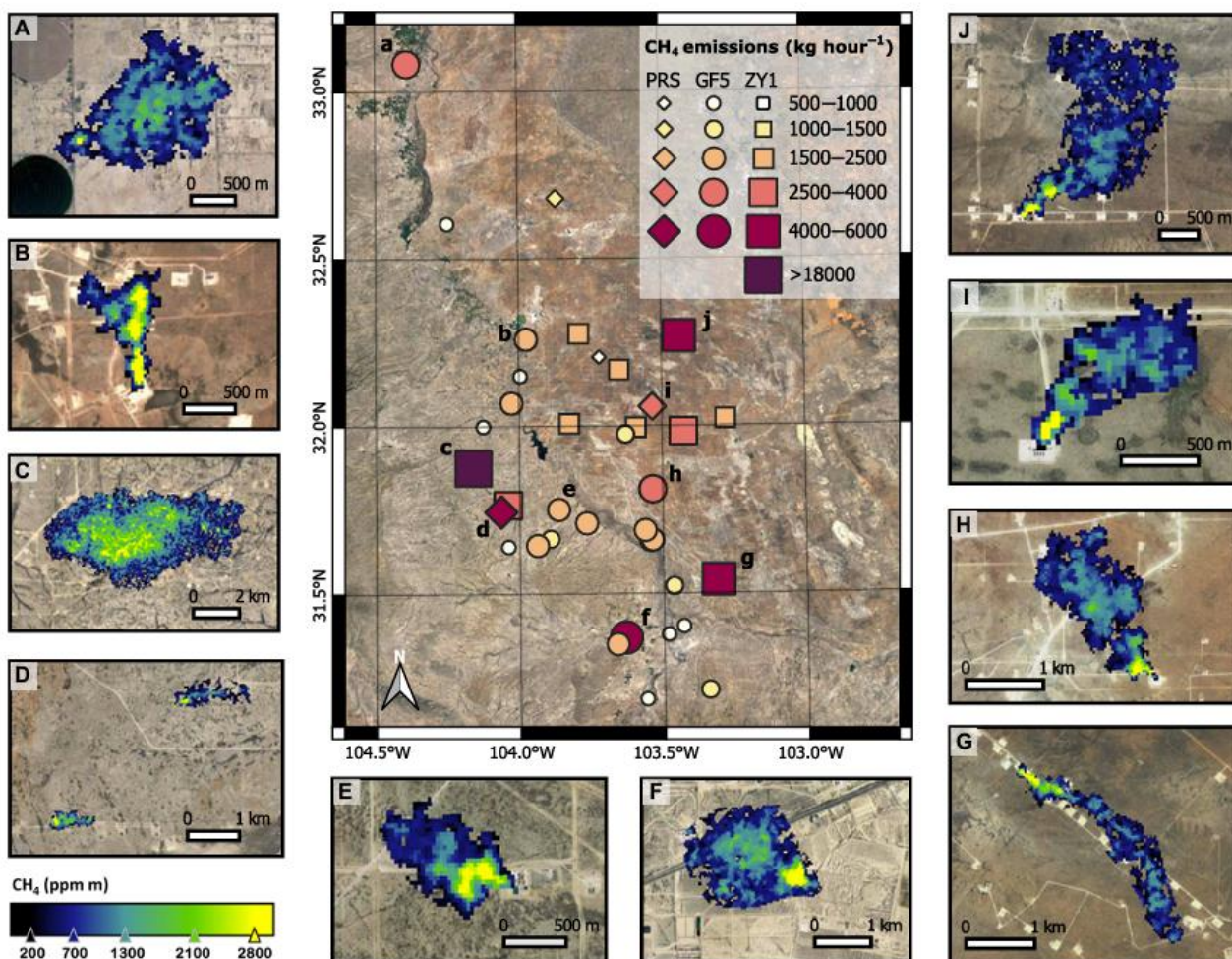


Figure 25 Extreme methane emissions detected in the Permian basin from satellite imaging spectroscopy data. A map with the identified methane plumes is shown in the central panel. Emissions are coded according to their flux rate and to the source of data (GF5-AHSI, GF5; ZY1-AHSI, ZY1; PRISMA, PRS). The small panels (A to J) around the main figure show examples of the detected plumes [48].

The researchers identified an unexpectedly large number of extreme point sources of methane emissions in the Permian basin, with 37 plumes having emission rates greater than 500 kg per hour, as seen in Figure 25. These sources accounted for a range between 31 and 53% of the estimated emissions in the sampled area. The study also revealed that new facilities are major emitters in the area, often due to inefficient flaring operations (20% of detections).

This work shows the potential application of using satellites for individual contributors, a prerequisite for mitigation.

IMPROVEMENTS BASED ON MACHINE LEARNING AND EXISTING DATASETS

This section focuses on how machine learning has been applied to improve the data quality from satellite measurements

A NEW MACHINE-LEARNING-BASED ANALYSIS FOR IMPROVING SATELLITE-RETRIEVED ATMOSPHERIC COMPOSITION DATA: OMI SO₂ AS AN EXAMPLE

Despite recent progress, satellite retrievals of anthropogenic SO₂ still suffer from relatively low signal-to-noise (SNR) ratios. In this study, we demonstrate a new machine learning data analysis method to improve the quality of satellite SO₂ products. This work presents a machine learning data analysis method aimed at improving the quality of satellite retrievals of anthropogenic sulfur dioxide (SO₂) emissions. The study utilizes slant column densities (SCDs) of SO₂ retrieved from the Ozone Monitoring Instrument (OMI) and calculates the ratio between the SCD and the root mean square (RMS) of the fitting residuals for each pixel. Pixels with low ratios, indicating presumably clean areas, are selected to build the training data with target SCDs set to zero. Pixels with higher ratios, indicating polluted areas, have target SCDs set to the original retrieved values. Neural networks (NNs) are then trained using predictors such as SCD/RMS ratios, solar zenith, viewing zenith and phase angles, scene reflectivity, and O₃ column amounts, as well as monthly mean SRRs. Two NNs are used for data analysis: one trained daily to produce analyzed SO₂ SCDs for polluted pixels each day, and the other trained monthly to produce analyzed SCDs for less polluted pixels for the entire month.

Test results show that the proposed method significantly reduces noise and artifacts over background regions, as seen in Figure 20. The monthly mean NN-analyzed and original SCDs generally agree within $\pm 15\%$ over polluted areas, indicating that the method retains SO₂ signals in the original retrievals, except for large volcanic eruptions. This is confirmed by using the NN-analyzed and original SCDs in a top-down emission algorithm to estimate annual SO₂ emissions for anthropogenic sources, which yield similar results. Alternative approaches using linear interpolation or a principal component analysis (PCA) - NN algorithm are explored, but they underestimate SO₂ over polluted areas.

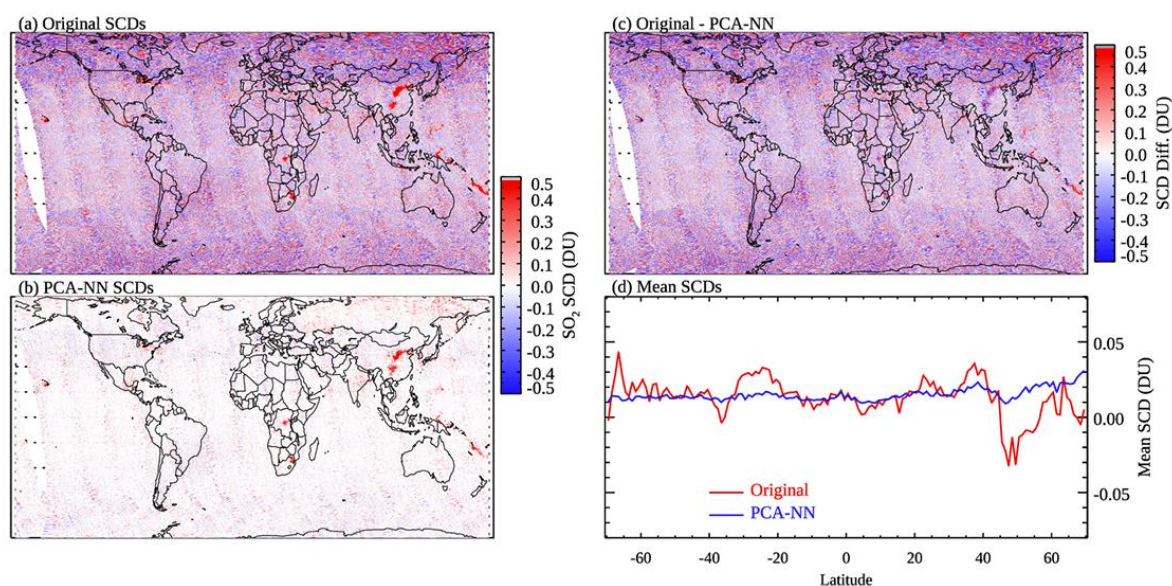


Figure 26 OMI SO₂ SCDs for 16 April 2005 retrieved using (a) the original PCA algorithm and (b) a PCA-NN algorithm, (c) the differences between the two retrievals, and (d) mean SO₂ SCDs for 1° latitude bands over relatively clean areas (monthly mean SRR < 3), calculated from (red) the original and (blue) PCA-NN retrievals [49].

Overall, the results demonstrate that the proposed data analysis method can improve the quality of existing OMI SO₂ retrievals, and it has the potential to be adapted for other sensors or species, enhancing the value of satellite data in air quality research and applications. This provides a use case for how machine learning can increase the data quality in satellite measurements.

SPECIFIC WORK ON THE DEVELOPMENT OF DATA-DRIVEN METHODS FOR IMPROVING METHANE EMISSION FROM SATELLITES

In this section, the focus will be on specific works that have been undertaken to develop data-driven methods for improving methane emission retrievals from satellites. These studies have utilized various approaches, including the use of neural networks, statistical models, and data assimilation techniques, to analyze satellite data and enhance the quality of methane emission estimates. The aim of these methods is to reduce noise, artifacts, and uncertainties in satellite retrievals, thereby improving our understanding of methane emissions from anthropogenic sources and supporting decision-making for mitigation strategies.

USING A DEEP NEURAL NETWORK TO DETECT METHANE POINT SOURCES AND 1 QUANTIFY EMISSIONS FROM PRISMA HYPERSPECTRAL SATELLITE IMAGES

In this work, a deep neural network was developed and designed to identify and quantify methane point source emissions from hyperspectral imagery from the PRecursores IperSpettrale della Missione Applicativa (PRISMA) satellite with 30-m spatial resolution. The network uses a combination of convolutional and fully connected layers to extract features from the satellite images and make predictions about the presence and quantity of methane emissions. The neural network was trained with simulated 26 synthetic methane plumes generated with the Large Eddy Simulation extension of the Weather Research and Forecasting 27 model (WRF-LES), which was embedded into PRISMA images [50].

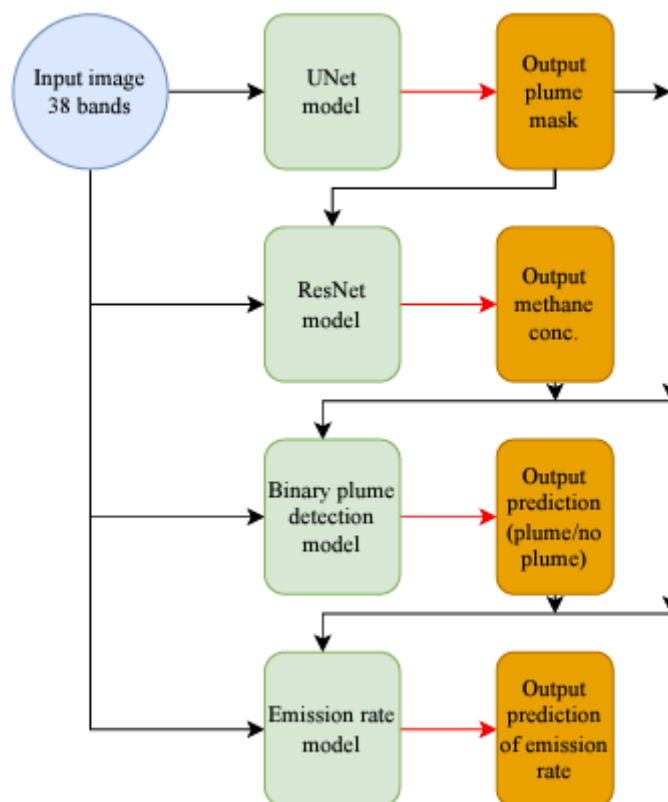


Figure 27 Structure of the neural networks used in this study. Green boxes indicate portions of the neural network, orange boxes indicate predictions made by each stage of the neural network. Black lines indicate flow of data into models, and red lines indicate predictions resulting from a model [50].

The model was then tested on 40 PRISMA scenes obtained during 2020–2022 in the Korpėje oil field, Turkmenistan, which is a well-studied area with frequent methane point source emissions plumes. The images were normalized in the same way that the training, test, and validation images were. 21 plumes were identified from 15 different scenes with predicted emission rates ranging from 1112–7615 kg per hour, shown in Figure 28.

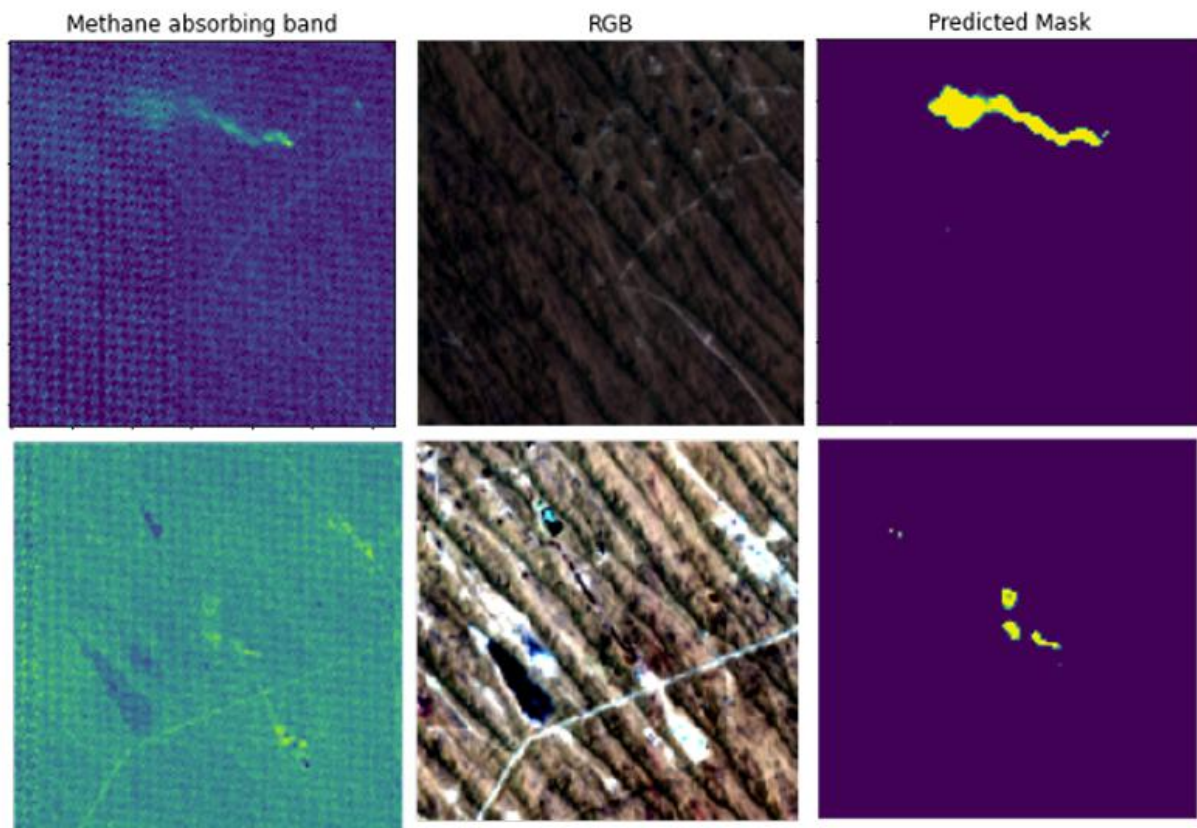


Figure 28 Images of plumes detected by the neural network in the Korpje oil field, Turkmenistan. Left panels depict methane 376 retrievals, middle panels depict the RGB of the image, and the right panel depicts the mask prediction by the neural network. The predicted emission rates are (top) 7615 and (bottom) 2370 kg hr⁻¹. RGB image courtesy of PRISMA © (Italian Space Agency) [50].

When compared to classical approaches like thresholding and clustering, deep neural networks have demonstrated greater success in identifying and quantifying methane point sources. The developed model exhibited impressive prediction capabilities, generating results for a 900 km² area using real PRISMA images in less than a minute. This remarkable speed and accuracy have the potential to significantly reduce the time and costs associated with mitigating anthropogenic methane emissions, which is a key factor for enabling a feasible method doing flare monitoring.

AUTOMATED DETECTION AND MONITORING OF METHANE SUPER-EMITTERS USING SATELLITE DATA

TROPOMI, a satellite instrument in orbit, has been providing daily global coverage of methane mixing ratios at a high resolution of up to 7x5.5 km², making it possible to detect methane super-emitters. However, the sheer volume of observations generated by TROPOMI, combined with the complexity of methane data, makes manual inspection impractical. To address this, this work has developed a two-step machine learning approach that utilizes a Convolutional Neural Network to detect plume-like structures in the methane data, followed by a Support Vector Classifier to differentiate between emission plumes and retrieval artifacts. These models are trained using data prior to 2021 and subsequently applied to all observations from 2021 [51].

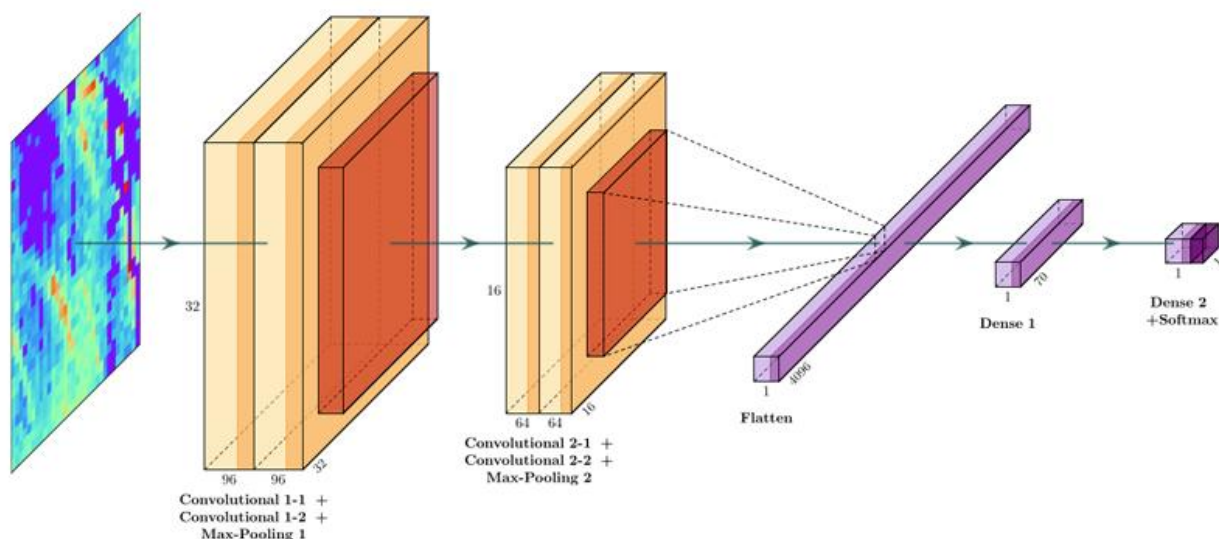


Figure 29 A schematic overview of the Convolutional Neural Network with a pre-processed 32x32 pixel TROPOMI methane scene (left) as input (Figure 1c). The CNN consists of two convolutional blocks (each with two Convolutional layers followed by a Max-Pooling layer) followed by two Dense (or Fully Connected) layers and an output node. Numerical values show input dimensions, layer dimensions and optimized hyperparameter values [51].

Furthermore, they employed a "tip-and-cue" approach for twelve (clusters of) TROPOMI detections in which one satellite instrument or data source "tips off" or identifies a potential emission event or anomaly, and then "cues" or directs another satellite instrument or data source to collect additional data or perform further analysis to confirm or pinpoint the exact source of the event. The "tip-and-cue" approach is used to identify the exact sources responsible for the plumes detected by TROPOMI, the satellite instrument used for methane monitoring. When TROPOMI detects a plume-like structure in the methane data, the system then uses high-resolution observations from other satellites such as GHGSat, PRISMA, and Sentinel-2 to collect additional data and perform detailed analysis to precisely identify the emission source associated with the detected plume. This approach helps in accurately attributing the emissions to specific facilities or locations, allowing for a more comprehensive understanding of the sources of methane emissions.

They utilize high-resolution observations from GHGSat, PRISMA, and Sentinel-2 satellites to pinpoint the exact sources responsible for these plumes. This has enabled us to detect and analyze both persistent and transient facility-level emissions underlying the TROPOMI detections, including emissions from landfills and fossil fuel exploitation facilities. In some cases, we have identified up to ten facilities contributing to a single TROPOMI detection.

This system detected 2974 plumes in 2021, with an estimated mean source rate of 44 t h⁻¹ and a 5-95th percentile range of 8-122 t h⁻¹. These emissions originate from 94 persistent emission clusters and numerous transient sources. Based on a comparison with bottom-up emission inventories, they have identified that most of the detected plumes are associated with urban areas/landfills (35%), followed by gas infrastructure (24%), oil infrastructure (21%), and coal mines (20%) [51] as seen in Figure 30.

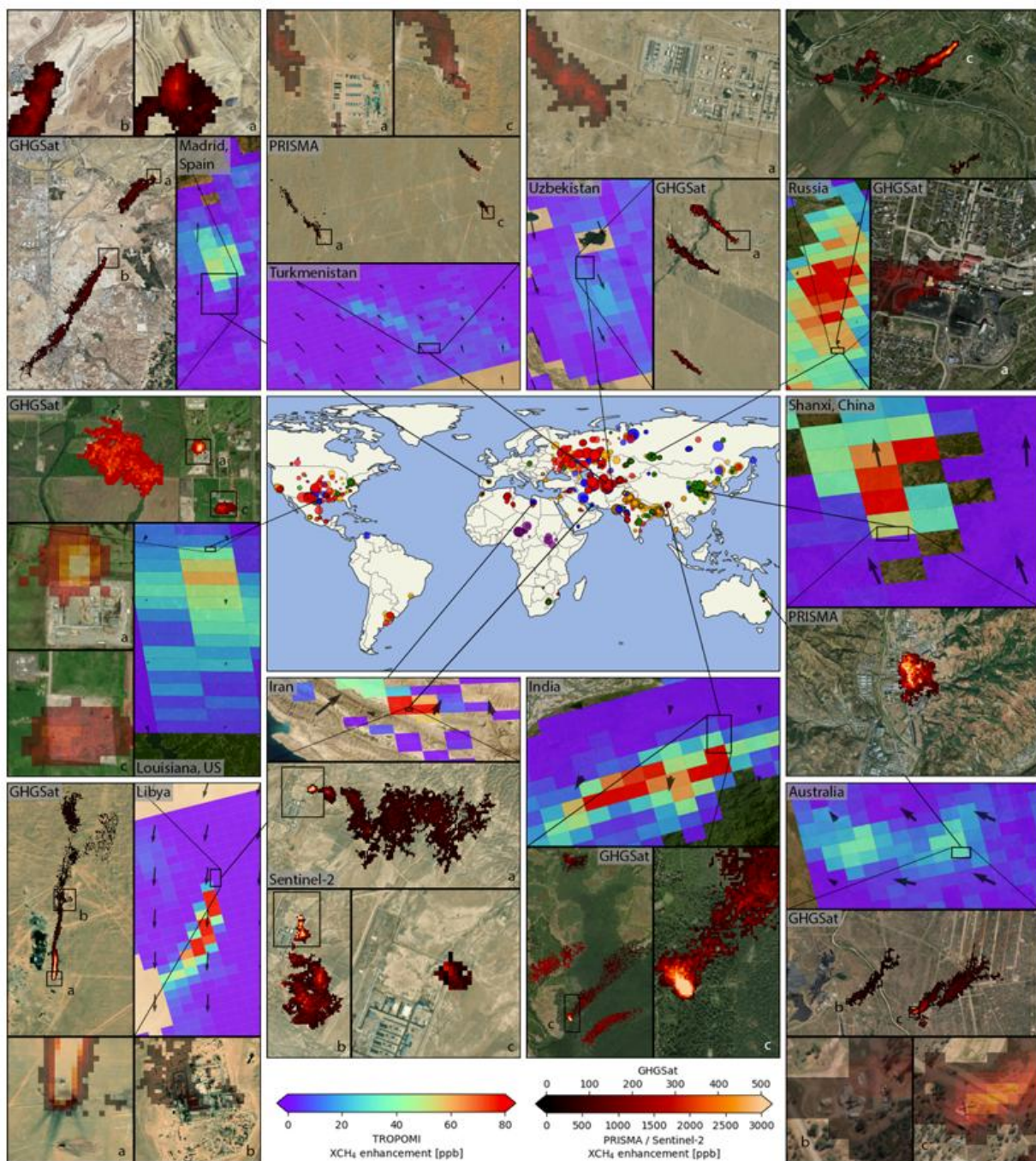


Figure 30 Plumes detected over 10 locations which were inspected with high-resolution instruments. Observations at the same location with different instruments are most often not on the same day [51].

This automated TROPOMI-based monitoring system, in combination with high-resolution satellite data, allows for the effective detection, precise identification, and continuous monitoring of methane super-emitters, which is crucial for mitigating their emissions, such as flaring.

SUMMARY

Based on this technical review, one can highlight that the current best practice for flare monitoring is using ultrasonic flow meters for estimating CO₂ emission with an assumed efficiency of 98% which is used for emissions reporting.

Employing satellite data for monitoring methane emissions is a novel approach, but there are some commercial offerings has started to be available, like from Maxar, which are not available to the public.

Further work is needed to ensure compliance with a methane tax scheme based on open datasets, such as those from the Copernicus program.

Based on the reviewed data, the performance matrix has been populated

Table 6 Performance matrix for the evaluated satellite platforms

Satellite	Data availability	Revisit time	Resolution
Sentinel-5	Suitable	Partly Suitable	Not suitable
Sentinel-2	Suitable	Partly Suitable	Suitable
Landsat-8	Partly Suitable	Partly Suitable	Partly Suitable
Worldview 8	Not suitable	Suitable	Partly Suitable
GHGSAT-METHANESAT	Not Suitable	Suitable	Partly Suitable
VIIRS	Partly Suitable	Partly Suitable	Not suitable

One can see, based on the reviewed data, that the VIIRS and Sentinel-2 satellite seems to be suitable for the data platform which this project can propose a method for enhancing the data quality based on a machine learning method to be more competitive or at least be the basis for a method that can be used in later satellite programs as a part of a wider knowledge transfer.

METHOD - DEVELOPMENT OF DATA -PREPROCESSING

In this chapter, the focus is on the development of a novel data pre-processing method aimed at estimating methane emissions from flaring activities in the Oil and Gas industry using satellite data.

The selection of suitable satellites for this work was based on the performance matrix as shown in Table 6. In it, different satellites are ranked based on Data availability, Revisit time, and resolution. Based upon the results in Table 6, Sentinel-5P and Sentinel-2 was selected as our desired satellite platforms, alongside data from VIIRS as a ground truth dataset with known flare locations.

Sentinel-2 was selected for the high resolution and the data availability, while Sentinel-5P was selected due to the data availability, and the partly suitable revisit frequency. However, it compensates for this by providing methane data, which is necessary in this project.

From the VIIRS satellite, a dataset derived from the VIIRS instrument will be used, containing a list of known flaresites, alongside other useful data.

The following subsections will review the available data sources from these satellites, before explaining a data pre-processing method, which results in a proof of concept.

REVIEW OF AVAILABLE DATA SOURCES

The basis for this pre-processing method is the application of Sentinel-2 and Sentinel-5P satellite data; therefore, in this section, a review of the Sentinel dataset will be performed to present the underlying characteristics of the data.

SENTINEL-2

Sentinel-2 (S2) is a wide-swath, high-resolution, multispectral imaging mission with a global 5-day revisit frequency. The S2 Multispectral Instrument (MSI) samples 13 spectral bands: visible and NIR at 10 meters, red edge, and SWIR at 20

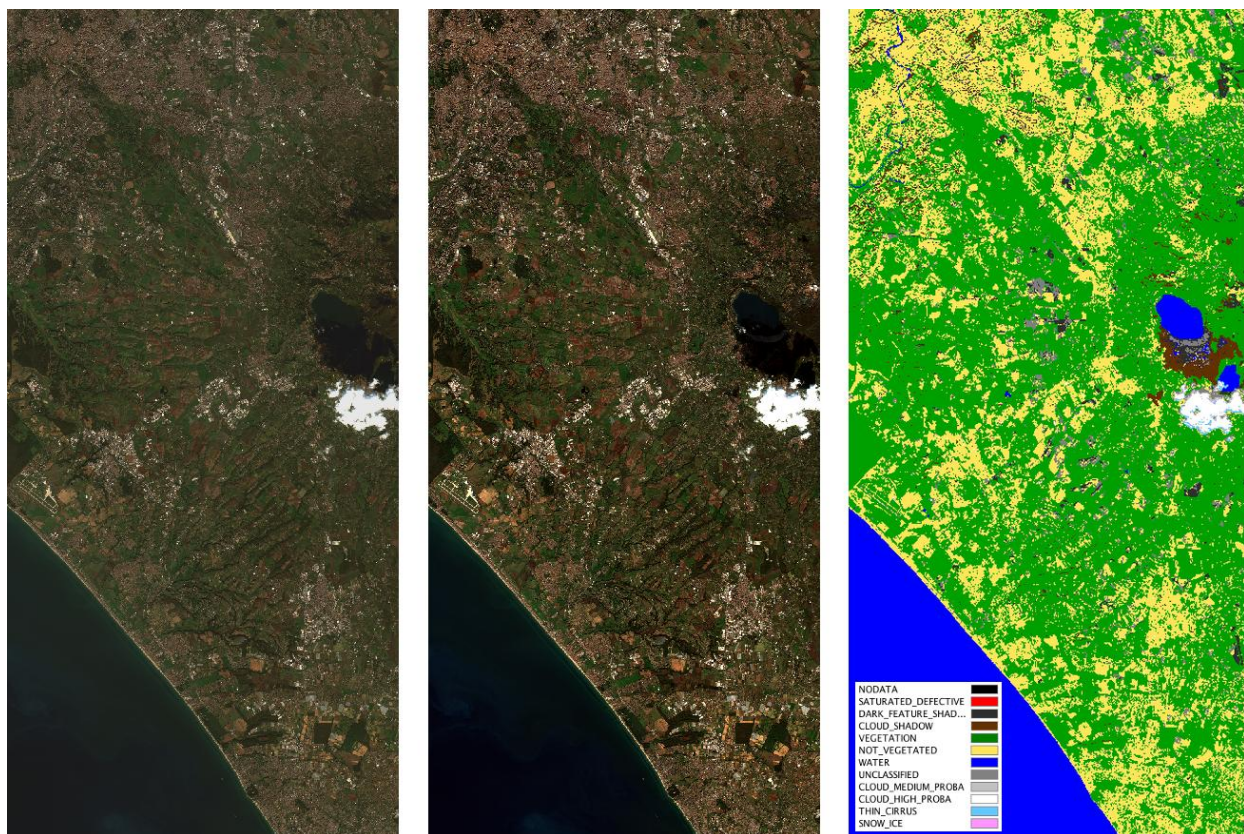
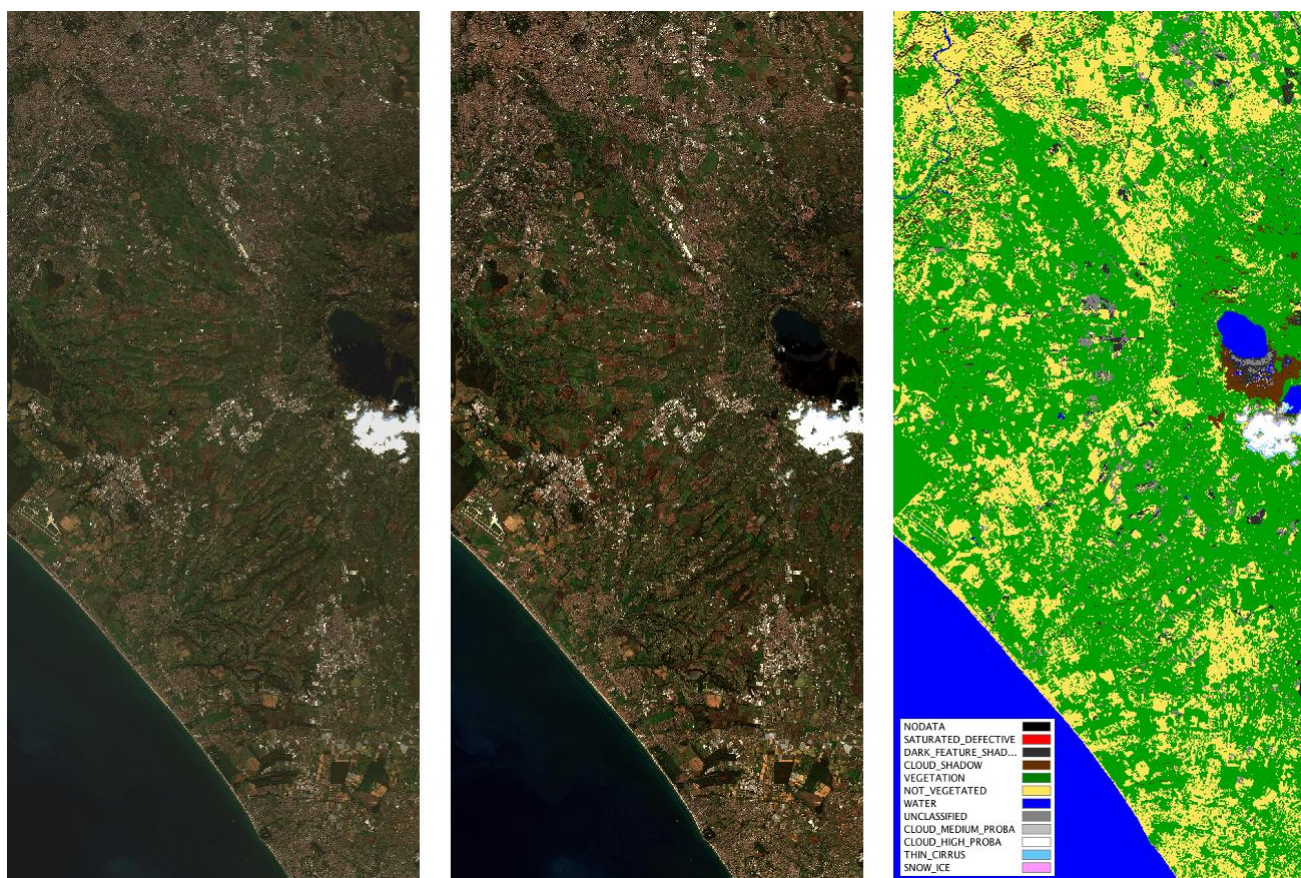


Figure 31 (1) Sentinel-2 Level-1C TOA reflectance input image, (2) the atmospherically corrected Level-2A surface reflectance image, (3) the output scene classification of the Level-1C product. Image taken from [39].

meters, and atmospheric bands at 60 meters spatial resolution. It provides data suitable for assessing the state and change of vegetation, soil, and water cover.



The Sentinel-2 Level-2A pre-Collection-1 product provides orthorectified Surface Reflectance (Bottom-Of-Atmosphere: BOA) with sub-pixel multispectral and multitemporal registration accuracy. Scene Classification (including Clouds), AOT (Aerosol Optical Thickness), and WV (Water Vapor) maps are included in the product. The main characteristics of the L2A products are listed in the following Table

Table 7 Sentinel-2 L2A Product Characteristics

Name	Level-2A
High-level Description	Surface reflectances in cartographic geometry for 12 spectral bands (10 m, 20 m, and 60 m depending on the wavelength; Cirrus band B10 is not included).
Algorithm used	Sen2Cor – Versions 2.10
Data Characteristics	UTM/WGS84 projection JPEG2000 image format 12-bit pixel values < 12 m at 95.5% confidence of Absolute Geolocation < 5 m at 95.5% confidence of Multitemporal Registration < 0.3 px at 99.7% confidence of Multispectral Registration
Additional Layers	Scene Classification Map at 60 m AOT Map Water Vapor Map
DEM used	Copernicus DEM at 30m
Auxiliary Data used	Ground Image Processing Parameters (GIPP) Digital Elevation Model (DEM) Global Reference Image (GRI) Copernicus Atmosphere Monitoring Service (CAMS) auxiliary parameters International Earth Rotation & Reference Systems Service (IERS) data Precise Orbit Determination (POD) data
Production & Distribution	Systematic generation and online distribution

Data Volume	800 MB (each 100x100 km ²)
Data Availability	Global since December 2018
Data Delivery	Available within 8 hours from sensing

The size of each tile is about 1.2 GB (compressed) and 4.8 GB (uncompressed), and the Sentinel-2 L2A products are available on the Copernicus Open Access Hub and through various data portals. These products also include a quality assessment band that offers valuable information regarding pixel quality [34]. This band helps identify factors that may affect the quality of the image, such as cloud and snow cover, water bodies, and shadows, allowing users to make informed decisions about the data's suitability for their specific applications.

SENTINEL-5 - TROPOMI (TROPOSPHERIC MONITORING INSTRUMENT).

As reviewed previously, the TROPOMI is a hyperspectral spectrometer with a nadir-viewing 108-degree Field-of-View push-broom grating, capable of measuring ultraviolet-visible (UV-VIS, 270nm to 495nm), near-infrared (NIR, 675nm to 775nm), and shortwave infrared (SWIR, 2305nm-2385nm) wavelengths.

The Sentinel-5P is designed to provide high-resolution measurements of atmospheric constituents such as ozone, NO₂, SO₂, CH₄, CO, formaldehyde, aerosols, and cloud properties. The retrieval algorithm used for the Sentinel-5P TROPOMI methane product is based on physics and utilizes the Oxygen-A band (760 nm) and absorption bands in the shortwave infrared spectrum. The algorithm includes a forward model of the optical properties of absorbing gases (oxygen, methane, water vapor, and carbon monoxide) and aerosols (size distribution, refractive index, and number concentration). The inversion is performed using the forward calculation, the measurement, and prior information, with cloud filtering being a critical step in the methane retrieval process.

The Sentinel-5P methane algorithm incorporates re-gridded cloud mask data from the Visible Infrared Imaging Radiometer Suite (VIIRS) for cloud filtering. When VIIRS cloud data are unavailable, additional filters based on Sentinel-5P/TROPOMI measurements and the FRESCO apparent surface pressure are applied. Other data filters in the retrieval algorithm include land-only pixels (excluding mountainous areas), spectrum intensity, solar zenith angle, and instrument zenith angle.

It is important to note that the Sentinel-5P/TROPOMI methane retrieval is intended for non-time-critical (NTC) data stream only. The main outputs of the retrieval process include the column-averaged dry air mixing ratio of methane, the random error, and the biased-corrected dry air methane fraction data based on the retrieved surface albedo.

Below is an image of the monthly average methane from January 2022, including the ocean glint retrieval.

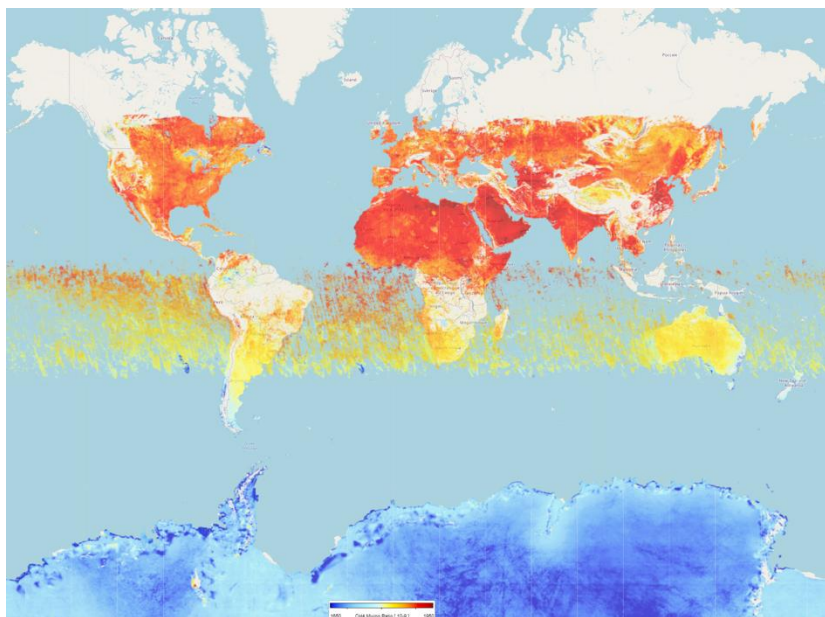
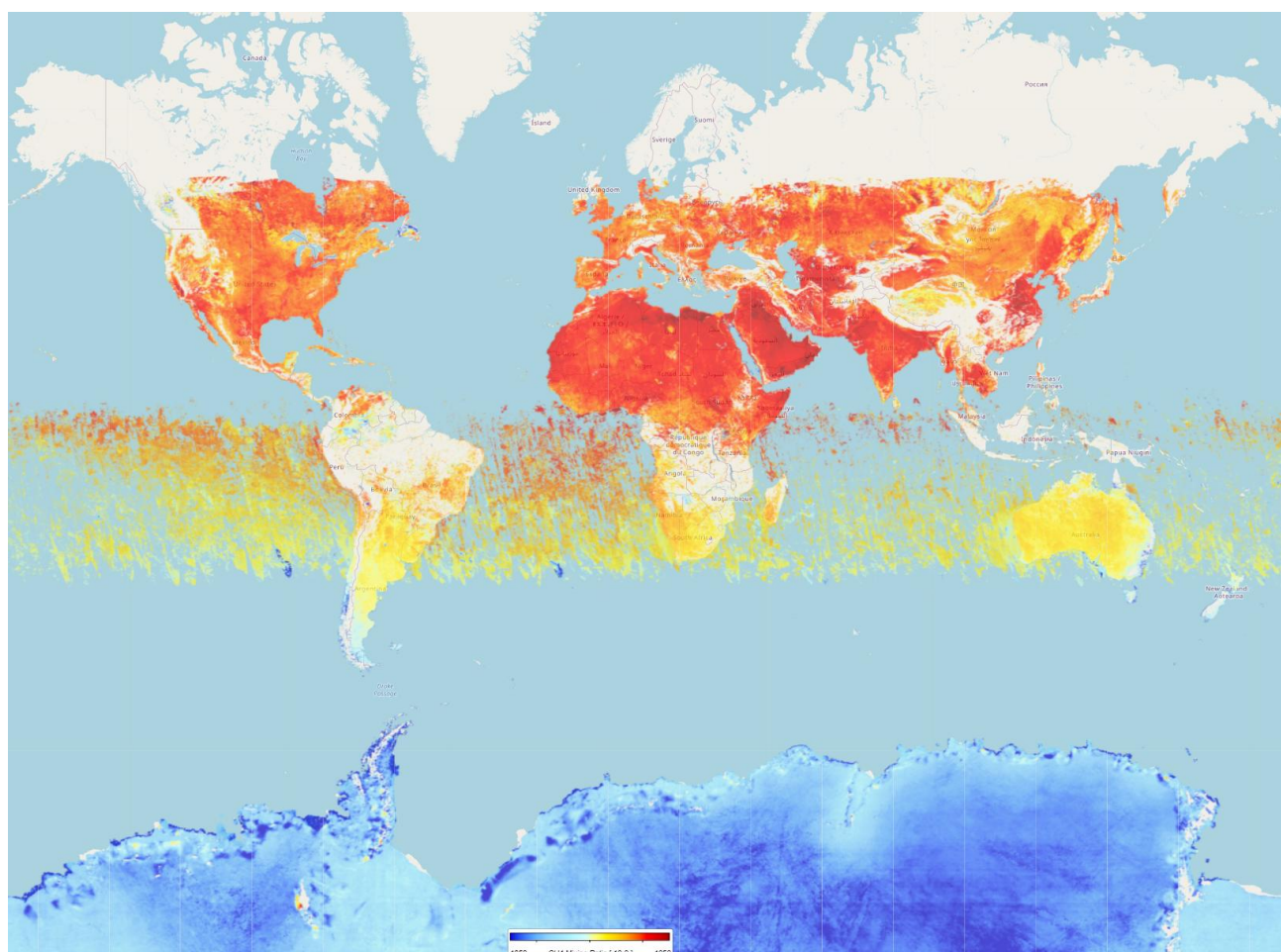


Figure 32 Monthly average methane from January 2022 including the ocean glint retrieval [40].



The methane data specifications for the Sentinel-5p satellite are as follows:

- **Spatial Resolution:** The TROPOMI instrument on Sentinel-5p provides a high spatial resolution of up to 3.5 km x 5.5 km, depending on the spectral band.
- **Spectral Coverage:** The TROPOMI instrument measures methane absorption in the shortwave infrared spectral range of approximately 1.6 to 1.7 micrometers.
- **Accuracy:** The TROPOMI instrument is designed to provide accurate measurements of methane concentrations in the atmosphere, with a target accuracy of better than 2% for column-averaged methane mixing ratios.
- **Temporal Coverage:** The Sentinel-5p satellite provides near-global coverage, revisiting the same location on Earth at least once per day. The TROPOMI instrument on Sentinel-5 measures methane concentrations daily, allowing for monitoring of short-term and long-term variations.
- **Data Format:** The methane data from Sentinel-5p is typically provided in netCDF (Network Common Data Form) format, a self-describing binary format widely used for storing and exchanging Earth observation data.

It is important to note that specific data specifications, including spatial resolution, accuracy, and data format, may vary depending on the version of the Sentinel-5 mission and the data product level (e.g., Level-1B, Level-2) being used. For the most up-to-date and accurate information, it is recommended to refer to the ESA's Sentinel-5p documentation or contact the relevant data providers for the latest data specifications.

Table 8 Sentinel-5 Methane Product Characteristics

Name	Sentinel-5P CH4: Methane
High-level Description	High-resolution imagery of methane concentrations.
Algorithm used	RemoTeC-S5P [52]
Data Characteristics	Swath-Based World-map projection TIFF, NetCDF, or 32-bit float data format.

	Radiometric accuracy of 1.6% in SWIR
Additional Layers	Geolocated total columns of ozone, sulfur dioxide, nitrogen dioxide, carbon monoxide, formaldehyde, and methane
DEM (Digital Elevation Model) used	None
Auxiliary Data used	ECMWF
Production & Distribution	Systematic generation and online distribution [32]
Data Volume	Ca 7.1 GB per day
Data Availability	Global since May 2018
Data Delivery	Available within 5 days of sensing

Further details of the Sentinel-5P CH4: Methane data specifications [53] can be found in the attachments.

GOOGLE EARTH ENGINE

The Google Earth Engine (GEE) API is an invaluable resource for obtaining and processing satellite imagery, including Sentinel-2 and Sentinel-5P data. This cloud-based platform allows researchers to efficiently access, analyze, and visualize large-scale geospatial datasets. By leveraging the power of the GEE API, this study can easily acquire satellite data and apply various processing techniques to prepare the data for further analysis.

One of the key advantages of using the GEE API is the availability of both Level 1 and Level 2 data. Level 1 data, also known as top-of-atmosphere (TOA) reflectance, provides raw, unprocessed images directly acquired from the satellite sensors. These images still contain atmospheric distortions and require additional processing steps to remove or reduce the impact of atmospheric effects. On the other hand, Level 2 data, also called bottom-of-atmosphere (BOA) reflectance, have already been processed to correct atmospheric influences. This pre-processing step, performed by the GEE API, dramatically enhances the usability and accuracy of the satellite data for various applications, including the study of methane emissions from flaring activities in the oil and gas industry.

The GEE API's ability to provide both Sentinel-2 and Sentinel-5P data in a pre-processed format allows researchers to focus on the analysis and interpretation of the data rather than spending time on pre-processing tasks. Moreover, by offering access to Level 1 and Level 2 data, the GEE API enables users to choose the most appropriate data type for their research objectives. In this work, using the GEE API for acquiring Sentinel-2 and Sentinel-5P data ensures a streamlined data acquisition process and facilitates the generation of a reliable dataset.

GEE provides a Python API, enabling developers and researchers to fetch and handle geospatial data from the GEE API programmatically.

VIIRS DATASET

The VIIRS derived dataset is obtained from the “Flaring Monitor” (FM) project, which is an open-source project that processes and relates NASA satellite sensor readings of heat signatures from natural gas flares to publicly disclosed ownership and reported flaring volumes to estimate the equivalent tons of CO2 emitted by companies in real-time. The FM repository contains raw and processed flare observation data for oil and gas-producing areas of the Lower 48 States of the US. The data contains processed VIIRS data, satellite-estimated Flaring Volumes, and state-reported Flaring Volumes both volumetrically and in their CO2 equivalents [54].

The FM project user interface provides an overview of the estimated monthly CO2 emission from different sites, as seen in the following Figure.

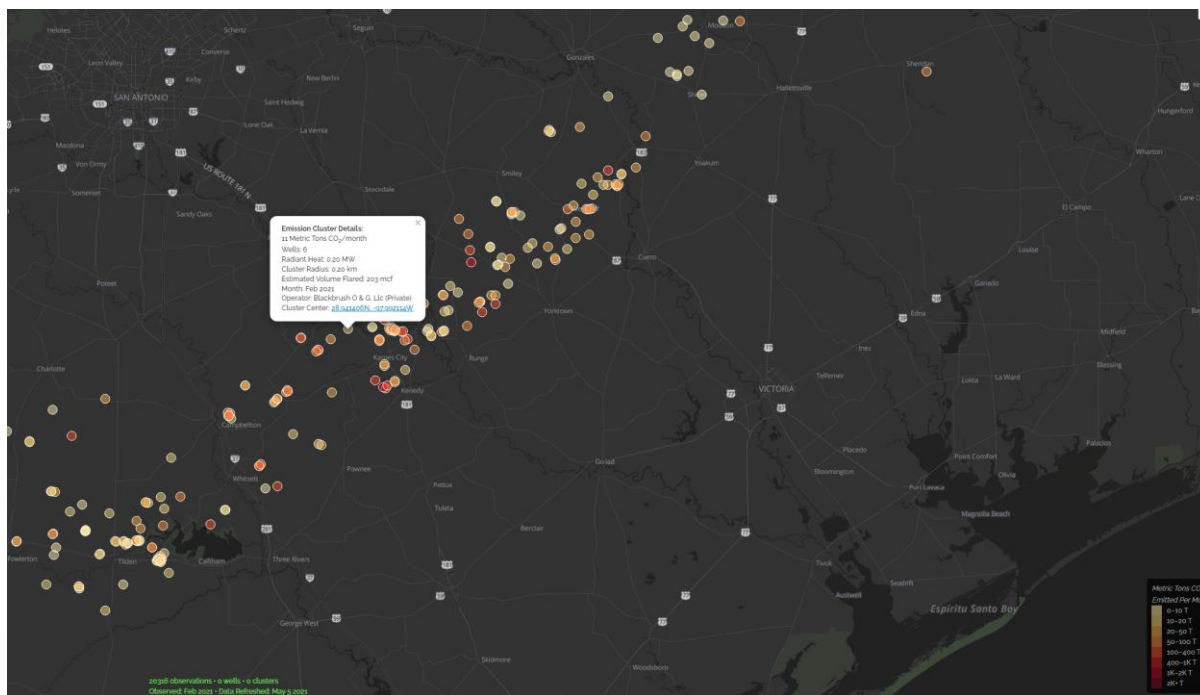


Figure 33 overview of the Flaring Monitoring user interface where the sites south of San Antonio has been used as a use case [55]

The dataset is divided up in two main sections:

1. Raw data
2. Processed data, which includes:
 - Processed VIIRS data
 - Satellite-estimated Flaring Volumes
 - State-reported Flaring Volumes (both volumetrically and in their CO2 equivalents)

This processed data is derived by combining VIIRS measurements with meticulously curated well meta and production data obtained from public regulatory authorities in New Mexico, Texas, North Dakota, and Colorado.

The datasets available is presented in two tables: Table 9 displays an overview of the raw datasets available from the VIIRS instrument, while Table 10 presents an overview of the processed datasets derived from the same instrument.

Both tables contribute to a comprehensive understanding of the dataset, allowing for a detailed exploration and interpretation of the VIIRS instrument data within the context of this study.

Table 9 Raw data from the VIIRS instrument

File name	Description
flare_prod.csv.gz	For a given API and month is well expected to be producing oil and/or gas
flare_wells.csv.gz	For a given API provides state, location, the most recent operator, and operator ticker if available.
lower_48_basins.csv.gz	EIA basin map

mapping_combined.csv.gz	Mapping of ticker to name
reported_flares.csv.gz	Monthly flaring as reported to state regulatory agencies. key_type and key provide a reference to the well or lease record in which this data was found.
states.csv.gz	State map
viirs.csv.gz	VIIRS radiant heat data for our four states (North Dakota, New Mexico, Texas, Colorado) region of interest

Table 10 Processed data from the VIIRS instrument.

File Name	Description
flaring_monitor_basin_stats.csv	Contains Flaring Monitor satellite-estimated Flaring Volumes and tons of equivalent CO ₂ emitted for the Permian, Eagle Ford, Denver Julesburg – Niobrara (Colorado only) and Bakken (North Dakota only). Data is available from 2019 onwards and is aggregated at the Basin level.
flaring_monitor_company_stats_satellite_modeled.csv	Flaring Monitor satellite-estimated Flaring Volumes and tons of equivalent CO ₂ emitted organized by Operator. The data is available for Operators in the states of Texas, New Mexico and North Dakota and is available from 2019 onwards.
flaring_monitor_company_stats_reported.csv	State reported flaring volumes and tons of equivalent CO ₂ emitted for Texas, New Mexico, Colorado, and North Dakota. The data is aggregated by Operator.
flaring_monitor_detailed_observations.csv	Unprocessed satellite radiant heat measurements from VIIRS in MW hours, associated estimated flared volume and tons of equivalent CO ₂ emitted by well cluster for operators in Texas, New Mexico, and North Dakota. Data is available from 2019 onwards.

The basis of the VIIRS dataset is provided in [56]. Table 11 provides an overview of the first nine rows of the VIIRS dataset, showcasing only the columns of interest. The dataset contains information such as the type of flaring activity, the month when the data was recorded, the sum of radiated heat (Sum RH), the estimated flare volume (in million cubic feet, MCF), the radius of the flaring area, the longitude (Long.) and latitude (Lat.) of the flare site, the number of wells, and the equivalent CO₂ released (in metric tons).

Table 11 An overview of the VIIRS dataset structure, taken from [57].

type	month	Sum RH	Est. Flare Vol (MCF)	Radius	Long.	Lat.	Wells	Equiv. CO ₂ released (Metric Tons)
mwc	2019-01-01	0.4835		0.2695	-102.7482	47.6163	11	
mwc	2019-01-01	6.6971	4289.8504	0.3569	-102.6316	47.7898	9	235.1
mwc	2019-01-01	27.665	20333.8305	0.3468	-103.1153	47.8755	4	1114.3
mwc	2019-01-01	15.5147	35202.2815	0.3409	-102.5636	47.5902	5	1929.1
mwc	2019-01-01	2.8938	1663.8903	0.2892	-103.1988	47.5734	1	91.2
....

DATA PREPROCESSING METHOD – DESCRIPTION

The data pre-processing method proposed in this work combines Sentinel-2 and/or VIIRS or Sentinel-5P satellite data sources to improve the accuracy and quality of methane emission estimation from flaring activities in the Oil and Gas industry. The Sentinel-2 satellite provides hyperspectral data (features) while Sentinel-5P offers mixing ratio CH₄ and dry air data (target). The pre-processing consists of a "data processing pipeline" with the following steps:

1. **Image Alignment:** The first step in the pre-processing pipeline is to align Sentinel-2 and VIIRS/Sentinel-5P satellite images to a common reference frame. This process corrects any spatial distortions or variations caused by the satellite imaging systems and ensures consistency in the subsequent analysis.
2. **Image Enhancement:** To improve image quality and visibility, various image enhancement techniques are applied, such as contrast stretching, histogram equalization, or gamma correction. These methods adjust the brightness and contrast of the images, enabling better visualization and more accurate feature extraction.
3. **Feature Extraction:** After segmenting the images, relevant features are extracted from the regions of interest. These features may include size, shape, and intensity of the flare, as well as contextual information like surrounding terrain and vegetation. The extracted features will be used as input for the machine learning model to predict methane emissions.

By implementing this data pre-processing method, this work aims to establish a robust and reliable foundation for the subsequent development and validation of a future machine learning workflow as proposed in Figure 35. The method will ensure that the model is trained on high-quality data that accurately represents the complex relationships between satellite-derived features and methane emissions from flaring activities in the Oil and Gas industry.

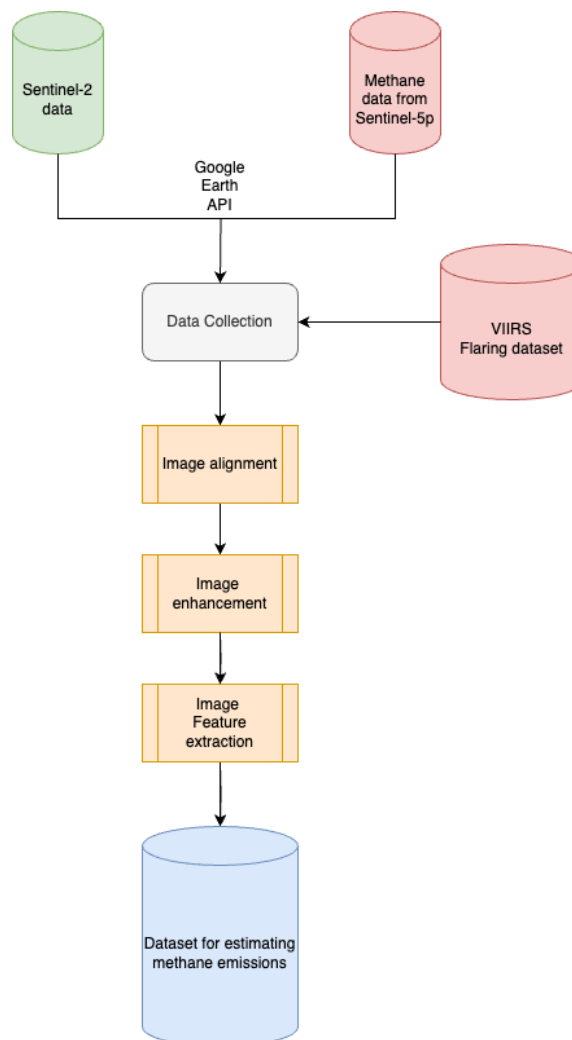


Figure 34 A high level overview of the dataset generation

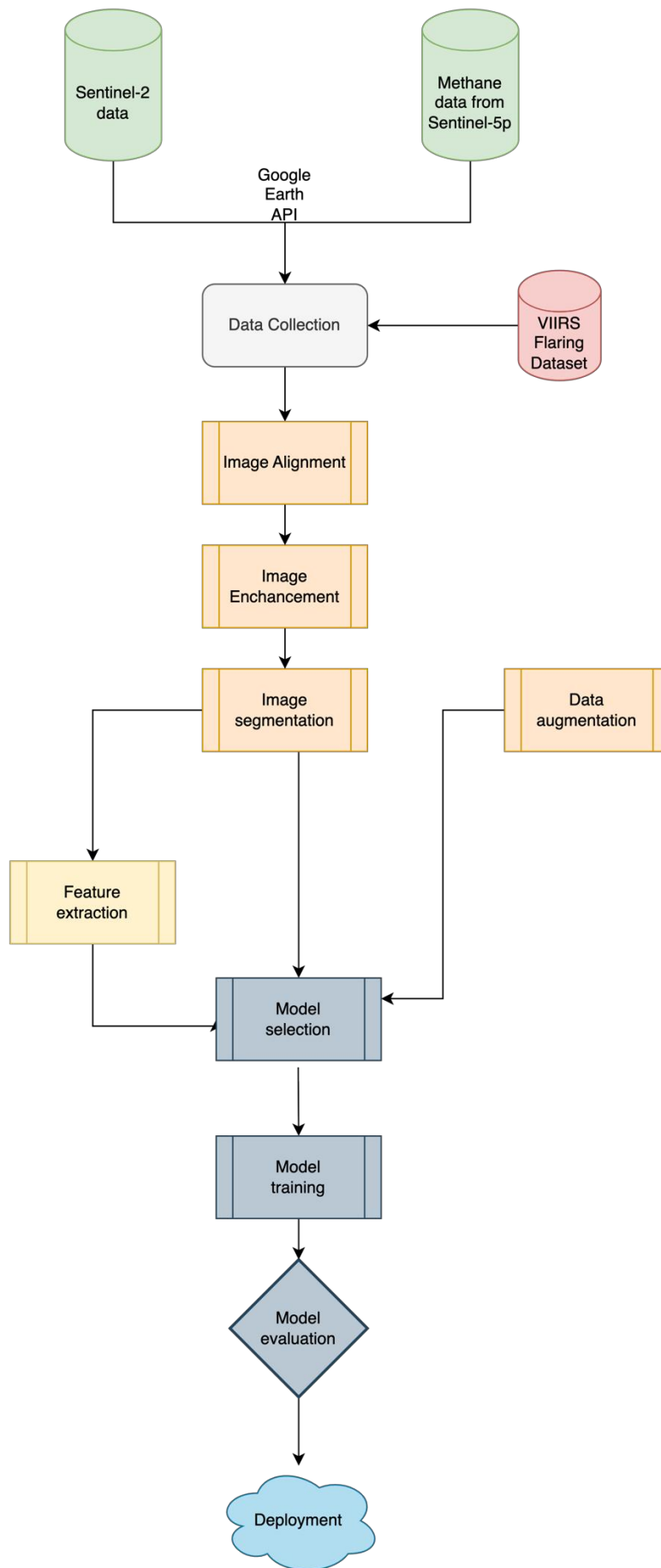


Figure 35 An overview of a proposed workflow for developing the pre-processing method based on machine learning.

DEVELOPMENT OF AN «PROOF OF CONCEPT»

In this section, we present a proof of concept that demonstrates the feasibility of our proposed method to enhance the accuracy and quality of methane emission estimation from flaring activities in the Oil and Gas industry using satellite data. The methodology is developed using Google Earth Engine API in Python, where satellite images from selected platforms are combined with the VIIRS dataset containing information about known flaring sites to generate a dataset of known flare sites.

DATA ACQUISITION AND PREPROCESSING

The satellite data used in this work are sourced from the Sentinel-2, and Sentinel-5 in addition VIIRS. In the data pre-processing stage, the *flaring_monitor_detailed_observations.csv* file is loaded, and the relevant columns are extracted: month, longitude, latitude, and *estimated_flare_volume_mcf*. This subset of data is used to initialize the Google Earth Engine and obtain the corresponding satellite data from Sentinel-2 and Sentinel-5 platforms.

Once the data has been acquired, the next step is to pre-process the data to prepare it for analysis. This may involve filtering the data based on specific criteria, such as cloud coverage, and applying various techniques to improve the quality of the data. In the case of Sentinel-2 data, cloud masking is applied by using the QA60 band in the Sentinel-2. This band helps identify clouds and their shadows, allowing us to remove or reduce their impact on the data. By applying cloud masking, we can ensure a more accurate representation of the underlying features, such as methane emissions from flaring activities, and improve the overall quality of the analysis.

In this work, a cloud coverage filter was applied to the satellite data to enhance its quality and reliability for further analysis. This filter ensures that only images with less than 20% cloud coverage are retained for the study. By limiting the data to images with minimal cloud interference, we can more effectively assess the features of interest, such as methane emissions from flaring activities, and draw more accurate conclusions from the analysis.

The side-by-side satellite images seen in Figure 36 shows the coast of Portugal and Spain, showcasing the significant impact of pre-processing on the clarity and usability of the data. In the image on the left, taken prior to pre-processing, the coastal regions appear white due to the presence of dense cloud cover. These clouds obstruct the view of the terrain and render the image less useful for analysis, as the true features of the landscape are hidden. However, in the image on the right, which has undergone pre-processing and cloud removal, the coastal regions are clearly visible, revealing the natural colours and features of the landscape. This stark contrast between the two images highlights the importance of pre-processing techniques, such as cloud masking, in ensuring that satellite data is accurate, reliable, and suitable for further analysis in various applications, including the study of methane emissions from flaring activities in the oil and gas industry.



Figure 36 Comparison of sentinel-2 data before and after preprocessing: The left image shows an image generated with the hyperspectral data with clouds, while the right image reveals a clear view of the coast after cloud removal. Both images are showing the B4, B3, and B2 bands.

The Sentinel-5P dataset in this work is provided through the Google Earth Engine (GEE) API, in an already pre-processed format, with the processing pipeline and methodology designed to be reproducible. This allows other researchers to replicate the results and conduct further investigations based on our proposed method.

An essential aspect of GEE's methodology is the filtering of the Sentinel-5P data based on a metric called *CH4_column_volume_mixing_ratio_dry_air_validity*, QA_{CH_4} , a continuous quality descriptor ranging from 0 (no data) to 100 (full quality data). This value serves as an indicator of the reliability and accuracy of the data, with higher values reflecting greater confidence in the measurements. We retain only data points with a validity value greater than 50, ensuring that the data used in our analysis meets a high-quality threshold, minimizing the impact of potential errors or inconsistencies [58].

In addition to the Sentinel-5P data, the VIIRS dataset utilized in this work is sourced from the Flaring Monitor's GitHub repository, specifically the *flaring_monitor_detailed_observations.csv file* [57]. This dataset contains unprocessed satellite radiant heat measurements from the Visible Infrared Imaging Radiometer Suite (VIIRS) in megawatt-hours (MWh), along with associated estimated flared volumes and tons of equivalent CO₂ emitted by well clusters for operators in Texas, New Mexico, and North Dakota. The data spans from 2019 onwards, providing valuable insights into the spatial and temporal trends of flaring activities in these regions.

In this work, the VIIRS dataset serves a crucial role in identifying known areas where flaring occurs, allowing for the targeted extraction of satellite data from Sentinel-2 and Sentinel-5 platforms. By combining the information from the VIIRS dataset with data from these other satellite sources, a comprehensive understanding of the environment and the impacts of flaring activities can be achieved. Although no specific preprocessing steps were applied to the VIIRS dataset, it still contributes valuable ground truth information to the analysis, helping guide the investigation and ensuring that the study focuses on relevant and significant areas of interest. This integrative approach, combining multiple data sources, enhances the overall reliability and robustness of the findings, enabling a more accurate assessment of methane emissions from flaring activities in the oil and gas industry.

IMAGE ALIGNMENT

In the proposed method, the image alignment process is significantly simplified due to the use of pre-processed Sentinel-2 Level-2 and Sentinel-5P Level-3 data obtained from Google Earth Engine. Both datasets have undergone various preprocessing steps, such as geometric and radiometric corrections, ensuring that the images are well-aligned and ready for further analysis.

Sentinel-2 Level-2 data includes atmospheric corrections and accurate geo-referencing, allowing for the precise overlay of images on maps or integration with other geospatial datasets. The Sentinel-5P Level-2 data on the other hand, is binned by time, rather than by latitude and longitude. To facilitate easy of usage with other geospatial datasets, GEE provides a Level-3 product of the data, where the original Level-2 data is converted, ensuring compatibility with other datasets.

This is done by applying equation (2), area-averaging values for a

The processing methodology used for the Sentinel-5P Level-3 data is openly available, enabling researchers to apply the same processing steps to their own Sentinel-5P Level-2 data provided by ESA if desired. This transparency promotes reproducibility and allows the scientific community to verify and build upon the findings of this study.

By using these pre-processed datasets, there is no need to align the images, due to the fact that Level-3 data is already aligned.

IMAGE ENHANCEMENT AND FEATURE EXTRACTION

In the Image Enhancement and feature extraction section of this study, the Normalized Burn Ratio (NBR) index is employed to improve the quality and usability of the satellite imagery acquired from Google Earth Engine (GEE). The main objective of this method is to highlight the areas impacted by gas flaring events, allowing for a more precise estimation of methane emissions from these activities.

To apply the NBR index within the GEE framework, the following high-level steps are applied on the processed Sentinel-2 data:

- **Band Selection:** The near-infrared (NIR) and shortwave infrared (SWIR) bands from the Sentinel-2 data are selected for further processing. These bands play a crucial role in determining the NBR index, as their distinct responses to vegetation and burned areas enable the detection of gas flaring impacts.
- **NBR Calculation:** The NBR index is calculated using equation (1). This calculation is performed using GEE's built-in mathematical functions, which allow for the efficient processing of large-scale geospatial datasets.

By applying the NBR index as an image enhancement method on the pre-processed, this study can effectively utilize satellite imagery to highlight effects on the environment from flaring activities. The enhanced images provide valuable insights into the environmental impacts of these events, contributing to a better understanding of the relationship between gas flaring and methane emissions. After the calculation of the NBR index we added the index values to the dataset as a separate layer.

RESULT – PRESENTATION OF THE POC

In this final section of the method chapter, a brief overview of the Proof of Concept (POC) is provided, showcasing the data-driven approach to enhance the accuracy and quality of methane emission estimation from flaring activities in the Oil and Gas industry. The methodology employs Sentinel-2, Sentinel-5P, and VIIRS satellite data, focusing on generating a comprehensive dataset that enables independent verification of methane emissions.

By providing a robust and comprehensive dataset, this methodology has the potential to contribute to better decision-making processes and the development of effective mitigation strategies related to flaring activities and methane emissions in the oil and gas industry. This data-driven approach can serve as a valuable tool for industry stakeholders, regulators, and researchers in their efforts to monitor and manage methane emissions from flaring activities.

A series of images is presented to illustrate the data transformation and improvement at each stage of the process:

Figure 37 shows the satellite data with atmospheric distortions and other noise that may impact the analysis, but with cloud cover removed. The removal of the clouds is necessary, since without it the images would just contain white pixels as seen in Figure 36.

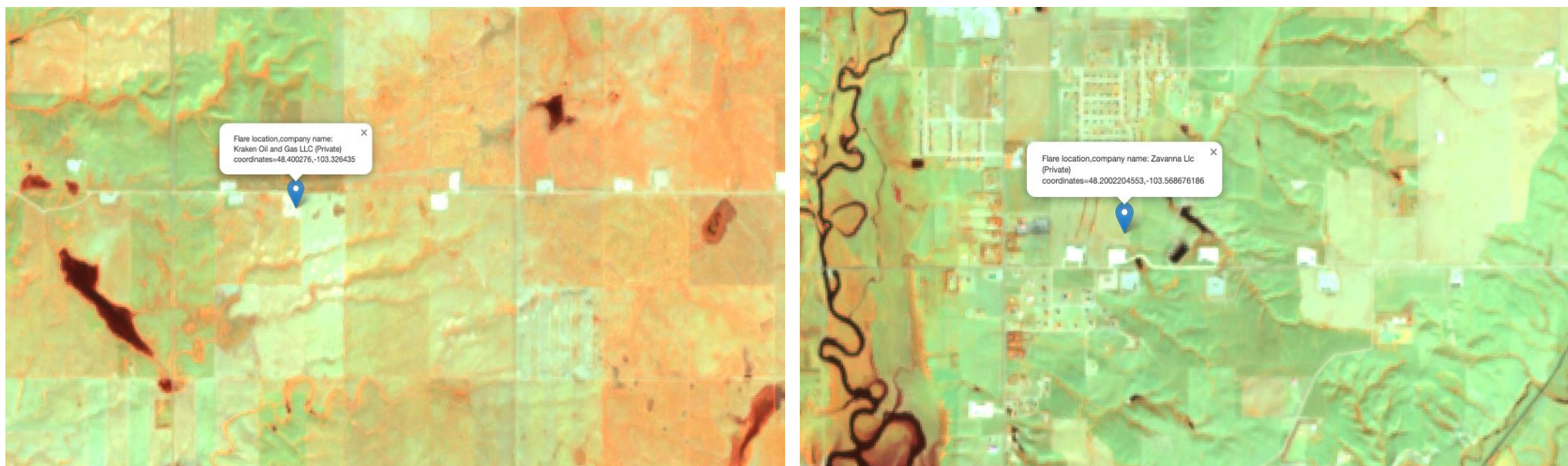


Figure 37 Sample images of the Raw Sentinel-2 data with the clouds removed, and with two areas of interest from the VIIRS dataset. On the left, we see a flare site from Kraken Oil and Gas LLC, and on the right, a flare site from Zavanna LLC.

Figure 38 displays the result of preprocessing techniques, such as cloud masking, applied to Sentinel-2 data, significantly improving the clarity and usability of the data for further analysis.

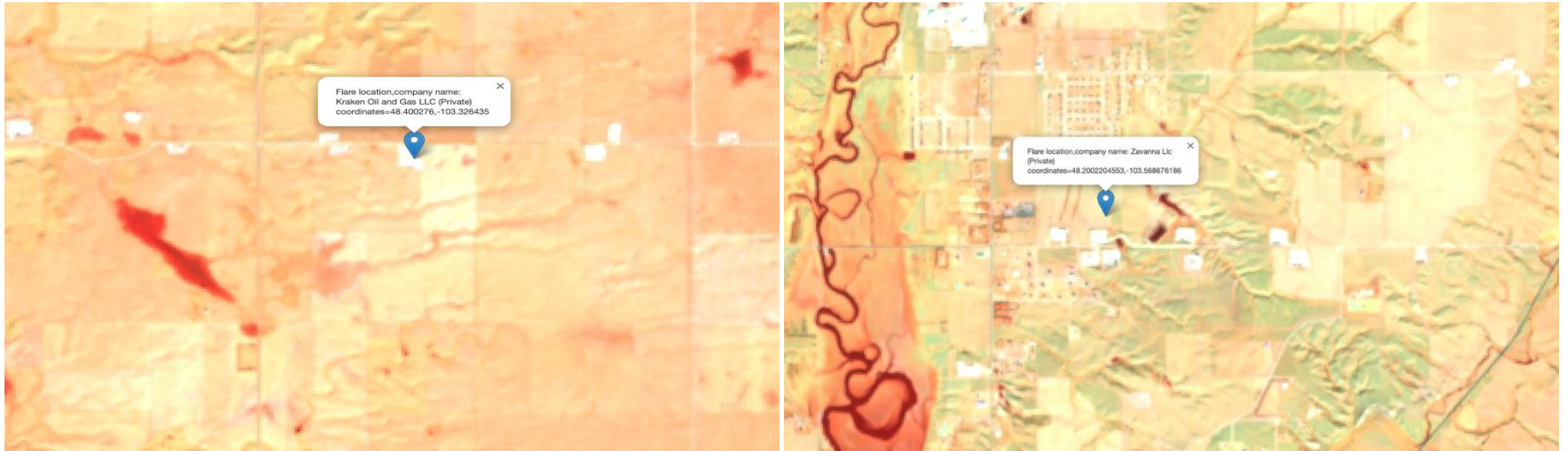


Figure 38 Sample image of the processed Sentinel-2 data, with two areas of interest from the VIIRS dataset. On the left, we see a flare site from Kraken Oil and Gas LLC, and on the right, a flare site from Zavanna LLC.

Figure 39 showcases the application of the NBR index as a separate layer to Sentinel-2 data. The NBR index aids in differentiating areas of potential interest related to flaring activities without asserting a direct relationship to methane emissions. This enhanced layer serves as an additional input for further analysis, providing context and support in the investigation of methane emission estimation from flaring activities while maintaining a neutral stance on the direct impact of the index on emission calculations.

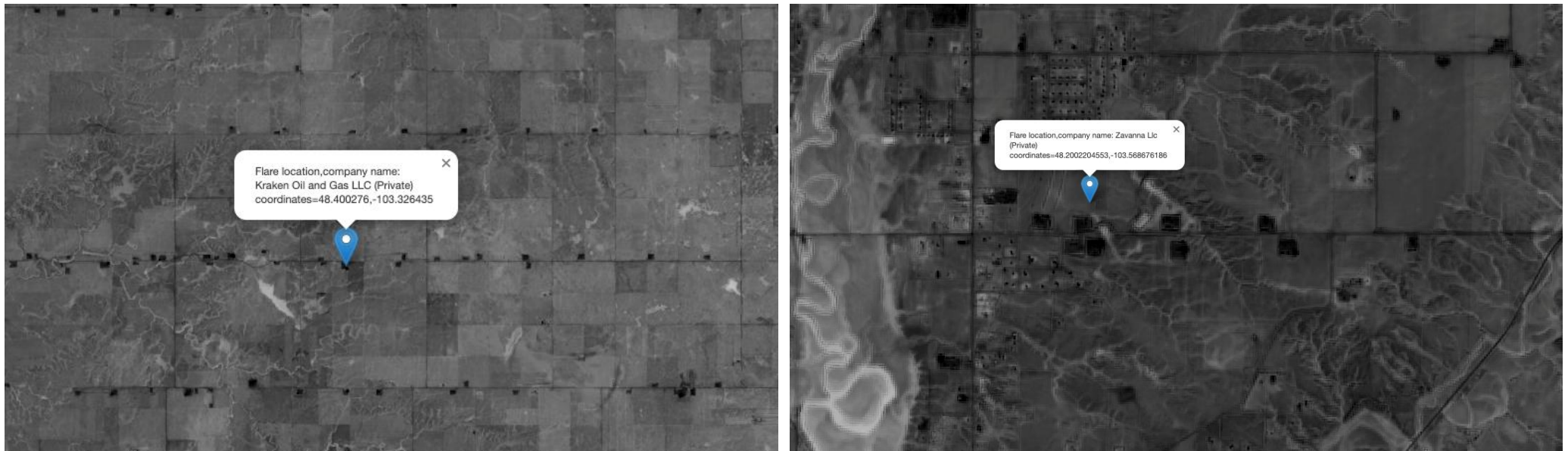


Figure 39 Sample image of the NBR index computed from the processed Sentinel-2 data, with two areas of interest from the VIIRS dataset. On the left, we see a flare site from Kraken Oil and Gas LLC, and on the right, a flare site from Zavanna LLC.

In the analysis of Sentinel-5P data, a customized color palette is utilized to effectively visualize the data, as depicted in Figure 40. The color palette represents a numerical value range of [1750, 1900], allowing for improved differentiation and interpretation of the data.



Figure 40 A representation of the color palette chosen for the Sentinel-5P data, corresponding to a numerical value range of [1750, 1900]

Figure 41 demonstrates the combined use of Sentinel-2 and Sentinel-5P data in the analysis of methane emissions from flaring activities. The left image of Figure 41 depicts the integration of Sentinel-2 multispectral data with Sentinel-5P methane concentration data. The right image of Figure 41 showcases the application of the Normalized Burn Ratio (NBR) index derived from Sentinel-2 data, alongside the Sentinel-5P data. Figure 40 highlights the methodological approach taken in this work, utilizing the combined datasets and processing techniques to effectively generate a comprehensive dataset for the analysis of methane emissions from flaring activities in the oil and gas industry.

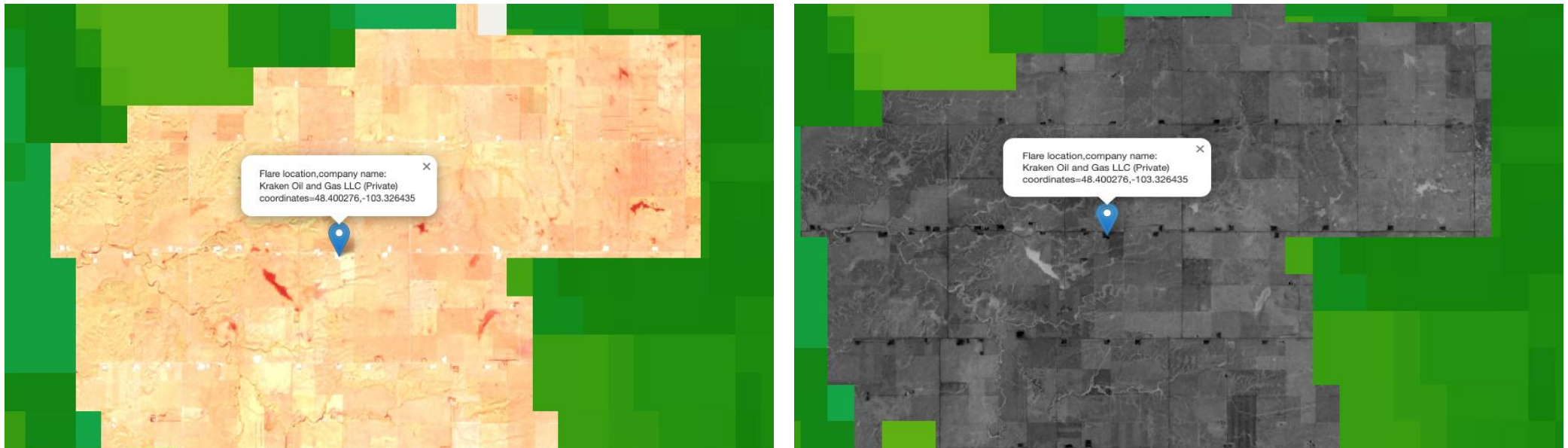


Figure 41 A visual representation of the integrated Sentinel-2 and Sentinel-5P data, showcasing the combined insights of both platforms. The left image presents the Sentinel-2 multispectral data combined with Sentinel-5P methane concentration data, while the right image displays the Normalized Burn Ratio (NBR) index derived from Sentinel-2 data merged with Sentinel-5P data. Both images show a flare site for Kraken Oil and Gas LLC

EVALUATION OF THE METHOD

In this section, we propose a test plan to further evaluate the effectiveness of the data-driven approach for enhancing the accuracy and quality of methane emission estimation from flaring activities in the oil and gas industry. The test plan aims to assess the performance of the proposed dataset and methodology in detecting and quantifying methane emissions from flaring activities, as well as its potential implications for the industry.

TESTPLAN FOR FUTURE TESTING

OBJECTIVES

The main objectives of the test plan are as follows:

- To evaluate the accuracy of the proposed dataset and methodology in estimating methane emissions from flaring activities.
- To assess the quality of the generated dataset in terms of its comprehensiveness, reliability, and usability.
- To examine the potential impact of the proposed approach on decision-making processes and mitigation strategies in the oil and gas industry.

TEST PLAN COMPONENTS

The test plan comprises several components designed to address the objectives outlined above. These components include:

GROUND TRUTH DATA COLLECTION

To assess the accuracy of the proposed dataset and methodology, ground truth data will be collected from a selection of flaring sites in the oil and gas industry. This data will serve as a reference point for comparing the satellite-derived estimates of methane emissions. Ground truth data collection methods may include direct measurements of methane emissions, such as gas chromatography or cavity ring-down spectroscopy, as well as other relevant data sources, such as industry reports and government records.

COMPARISON WITH EXISTING METHANE EMISSION ESTIMATION TECHNIQUES

To evaluate the proposed dataset and methodology, a comparison will be made with existing techniques for estimating methane emissions from flaring activities. This may include comparing the results obtained from the proposed approach with those derived from alternative satellite data sources, models, and ground-based measurements. The comparison will provide insights into the relative performance of the proposed methodology, as well as areas of potential improvement.

DATA QUALITY ASSESSMENT

The quality of the generated dataset will be assessed in terms of its comprehensiveness, reliability, and usability. This will involve examining various aspects of the dataset, such as the spatial and temporal resolution, coverage, and consistency of the data. Additionally, the quality of the pre-processing steps, such as cloud masking, atmospheric correction, and image enhancement techniques, will be evaluated.

ANALYSIS OF POTENTIAL IMPLICATIONS FOR THE OIL AND GAS INDUSTRY

The test plan will also examine the potential implications of the proposed approach for decision-making processes and mitigation strategies in the oil and gas industry. This may involve conducting interviews with industry stakeholders, regulators, and researchers to understand their perspectives on the utility and relevance of the generated dataset. Furthermore, a review of current policies, regulations, and industry practices related to flaring activities and methane emissions will be conducted to assess the potential impact of the proposed approach.

TESTPLAN

For the evaluation of the proposed dataset and methodology's accuracy, the following steps could be followed:

1. Procure ground truth data from chosen flaring sites via direct measurement methodologies such as gas chromatography or cavity ring-down spectroscopy.
2. Supplement the primary data with additional data sourced from industry reports and governmental records.
3. Conduct a comparative analysis of the results derived from the proposed method against those obtained from other satellite data sources, models, and ground-based measurements.

To assess the quality of the generated dataset, one could:

1. Undertake a comprehensive review of the dataset, scrutinizing aspects such as its spatial and temporal resolution, coverage, and consistency.
2. Conduct a critical evaluation of the pre-processing steps, including but not limited to, cloud masking, atmospheric correction, and image enhancement techniques.

To evaluate the potential ramifications of the proposed approach on the industry, the following steps are recommended:

1. Conduct structured interviews with industry stakeholders, regulators, and researchers to solicit their expert perspectives on the proposed approach.
2. Review current policies, regulations, and industry practices related to flaring activities and methane emissions.
3. Conduct a comparative analysis between the potential impacts of the proposed approach and current industry practices to ascertain the feasibility and implications of the proposed method.

DISCUSSION

This section discusses the main findings and implications of the study, focusing on the data-driven approach developed to enhance the accuracy and quality of methane emission estimation from flaring activities in the oil and gas industry. The methodology employs Sentinel-2, Sentinel-5P, and VIIRS satellite data to generate a comprehensive dataset, enabling independent verification of methane emissions. This data-driven approach has the potential to contribute to better decision-making processes and the development of effective mitigation strategies related to flaring activities and methane emissions in the oil and gas industry.

COVERAGE DIFFERENCES BETWEEN SENTINEL-2 AND SENTINEL-5P

One of the key observations from the results is the difference in coverage between Sentinel-2 and Sentinel-5P data, which can be observed in Figure 40. While Sentinel-2 data provides comprehensive coverage of the areas of interest, Sentinel-5P data appears to be patchier, with less extensive coverage. This discrepancy may affect the accuracy of methane emission estimations in areas where Sentinel-5P data is limited or unavailable. Further research should focus on addressing this issue, potentially by exploring alternative satellite data sources or developing techniques to fill in the gaps in the Sentinel-5P data.

PREPROCESSING TECHNIQUES AND THEIR EFFECTIVENESS

The preprocessing techniques, such as cloud masking, applied to Sentinel-2 data significantly improve the clarity and usability of the data for further analysis. However, it is important to recognize that the effectiveness of these techniques may be dependent on the specific characteristics of the flaring sites being analyzed, as well as the quality of the input data. Future work should investigate the potential impact of these factors on the overall accuracy of the approach and explore other preprocessing techniques that may enhance the quality of the data.

APPLICATION OF THE NBR INDEX

The application of the Normalized Burn Ratio (NBR) index as a separate layer to Sentinel-2 data aids in differentiating areas of potential interest related to flaring activities without asserting a direct relationship to methane emissions. This enhanced layer serves as an additional input for further analysis, providing context and support in the investigation of methane emission estimation from flaring activities. However, the NBR index is an indirect measure of methane emissions, and its relationship with actual emissions may not be straightforward. Future research should seek to better understand this relationship and potentially explore alternative indices or methods for estimating methane emissions.

INTEGRATION OF SENTINEL-2 AND SENTINEL-5P DATA

The integration of Sentinel-2 multispectral data with Sentinel-5P methane concentration data, as demonstrated in Figure 40, highlights the methodological approach taken in this work. Despite the patchy nature of Sentinel-5P data, the combined use of datasets and processing techniques effectively generates a comprehensive dataset for the analysis of methane emissions from flaring activities in the oil and gas industry. However, it is crucial to acknowledge the limitations of the current approach and explore ways to improve the methodology and enhance the quality of the dataset.

RECOMMENDATIONS AND FURTHER WORK

This section discusses recommendations and further work based on the methodology developed in this work, which aimed to enhance the accuracy and quality of methane emission estimation from flaring activities in the oil and gas industry using a data-driven approach. The main approach focuses on utilizing satellite data from selected platforms, such as Sentinel-2, Sentinel-5P, and VIIRS, to generate a dataset suitable for.

Throughout this study, we have developed and applied a comprehensive methodology, encompassing data acquisition, preprocessing, image enhancement, and analysis. This methodology is designed to ensure the highest quality and reliability of the data used in our analysis, paving the way for more accurate and effective methane emission estimation.

Data acquisition is a critical step in our approach, and we leverage the capabilities of the Google Earth Engine (GEE) API to access Sentinel-2 and Sentinel-5P data efficiently. In addition, we utilize the VIIRS dataset to identify known areas of flaring activity and extract satellite data from these locations.

The processing stage involves a series of techniques aimed at improved the quality of the satellite data. For the Sentinel-2 data, we apply cloud masking using the QA60 band to minimize the impact of clouds on the dataset. We also filter the data ensuring that only high-quality data points are retained in the dataset.

Image enhancement techniques, such as the Normalized Burn Ratio (NBR) index, are applied to the processed Sentinel-2 data to further refine the information related to methane emissions from flaring activities. The NBR index highlights areas of potential methane emissions and is added as a separate layer to the data, enabling end-users of the dataset to remove this layer if they want to do so.

Finally, our approach integrates the processed Sentinel-2 and Sentinel-5P data, offering a comprehensive view of methane emissions from flaring activities. This integrated dataset enables more accurate estimation and detection of future emissions, contributing to the overall effectiveness of our data-driven approach.

FURTHER WORK

While the methodology developed in this thesis has shown promising results, there are several areas that could be addressed in future work to further enhance the quality and accuracy of methane emission estimates.

IMAGE SEGMENTATION

In the context of satellite imagery, image segmentation is a crucial step in the process of identifying and extracting regions of interest. The application of more advanced image segmentation techniques could enhance the detection of flaring sites. Future work could explore different algorithms and machine learning techniques, such as convolutional neural networks (CNNs), for this purpose.

FEATURE EXTRACTION

Feature extraction is another critical step in understanding the data and identifying flaring activities. Currently, this study uses the Normalized Burn Ratio (NBR) index as a feature. However, there are numerous other features that could be extracted from the data that may enhance the accuracy of methane emission estimates. Future research could explore the application of feature selection techniques to determine the most relevant features for this task.

DATA AUGMENTATION

The data augmentation techniques employed in this study were relatively simple. Given the advancements in this field, more sophisticated data augmentation methods could be explored in future work. Techniques such as geometric transformations, principal component analysis-based data augmentation, or generative adversarial networks (GANs)

could potentially provide more diverse and extensive data for training models, improving their performance and robustness.

EXPLORATION OF OTHER SATELLITE DATA SOURCES

This study was limited to Sentinel-2, Sentinel-5P, and VIIRS datasets. Future work could investigate using other satellite data sources that may provide additional or more detailed information relevant to methane emissions from flaring activities.

Other satellite data sources could be existing satellites like:

- Landsat 8 and 9

Or future satellites like:

- CHIME: Copernicus Hyperspectral Imaging Mission for the Environment
- CO2M: Copernicus Anthropogenic Carbon Dioxide Monitoring

INTEGRATION OF OTHER GROUND TRUTH DATA

Finally, including ground truth data could significantly enhance the validation of the developed method. Future studies could access such data with oil and gas industry partners to include a more diverse set of areas outside of North America

BIBLIOGRAPHY

- [1] U. N. E. Programme, «Methane emissions are driving climate change. Here's how to reduce them,» *United Nations Environment Programme. August*, vol. 20, p. 2021, 2021.
- [2] R. Lindsey og M. Scott, «After 2000-era plateau, global methane levels hitting new highs,» 2017.
- [3] European Commission, *Methane emissions - Energy*.
- [4] California State Senate, *Senate Bill 1137*, 2021.
- [5] A. E. Ramsden, A. L. Ganesan, L. M. Western, M. Rigby, A. J. Manning, A. Foulds, J. L. France, P. Barker, P. Levy, D. Say og others, «Quantifying fossil fuel methane emissions using observations of atmospheric ethane and an uncertain emission ratio,» *Atmospheric Chemistry and Physics*, vol. 22, p. 3911–3929, 2022.
- [6] European Commission, *Methane Emissions*, 2022.
- [7] N. E. Observatory, *Mapping Methane Emissions from Fossil Fuel Exploitation*, 2021.
- [8] S. M. Khodayee, F. Chiacchio og Y. Papadopoulos, «A Novel Approach Based on Stochastic Hybrid Fault Tree to Compare Alternative Flare Gas Recovery Systems,» *IEEE Access*, vol. 9, pp. 51029-51049, 2021.
- [9] S. E. Birkeland, G. Veire, C. Holm, H. J. Samuelsen, J. A. Torgersen, H. Nordang, V. Lossius og N. Mjølnørød, «Electrification and Other Measures to Minimize Carbon Emissions from the Johan Sverdrup Field,» 2020.
- [10] J. Elkind, E. M. Blanton, H. D. van der Gon, R. Kleinberg og A. Leemhuis, «Nowhere to Hide: Implications for Policy, Industry, and Finance of Satellite-Based Methane Detection,» *Columbia School of Public and International Affairs Center on Global Energy Policy*, 2020.
- [11] C. Brunetti, B. Dennis, D. Gates, D. Hancock, D. Ignell, E. K. Kiser, G. Kotta, A. Kovner, R. J. Rosen og N. K. Tabor, «Climate change and financial stability,» 2021.
- [12] I. E. Agency, *Inflation Reduction Act 2022: Sec. 60113 and Sec. 50263 on Methane Emissions Reductions*, 2022.
]
- [13] J. L. Ramseur, «Inflation reduction act methane emissions charge: In brief,» *Congressional Research Service,(R47206)*, 2022.
- [14] European Space Agency, *Types of orbits*, 2020.
]
- [15] T. G. Roberts, *Earth Orbit 101*, accessed April 8, 2023.
]

- [16 Anka Geo, *Satellite Imagery Technology*, accessed April 14, 2023.
]
- [17 European Space Agency, *About Payload Systems*, accessed 2023.
]
- [18 University of Hawaii, *3.4 Payload Design*, accessed 2023.
]
- [19 Xenics Photonics group , «Importance of Methane Detection and the use of Infrared to detect Methane Leak,» 02
] 02 2023. [Internett]. Available: <https://www.xenics.com/importance-of-methane-detection-and-the-use-of-infrared-to-detect-methane-leak/>.
- [20 Mathworks , «What Is Reinforcement Learning?,» 14 04 2023. [Internett]. Available:
] <https://se.mathworks.com/discovery/reinforcement-learning.html>.
- [21 S. Raschka og V. Mirjalili, Python machine learning: Machine learning and deep learning with Python, scikit-learn,
] and TensorFlow 2, Packt Publishing Ltd, 2019.
- [22 IBM, «What is a neural network?,» [Internett]. Available: <https://www.ibm.com/topics/neural-networks>. [Funnet
] 14 04 2023].
- [23 A. Ng, J. Ngiam, C. Foo, Y. Mai og C. Suen, «Multi-Layer Neural Network,» *UFLDL Tutorial*, 2013.
]
- [24 e. a. Alon Dadon, «Examination of spaceborne imaging spectroscopy datautility for stratigraphic and lithologic
] mapping,» *Society of Photo-Optical Instrumentation Engineers v0l 5*, p. 3507, 03 2011.
- [25 e. setyawan, 12 08 2019. [Internett]. Available: [https://www.intermap.com/blog/satellite-imagery-resolution-vs.-
\] accuracy](https://www.intermap.com/blog/satellite-imagery-resolution-vs.-accuracy). [Funnet 12 02 2023].
- [26 Maxar Technologies , «Spatial Resolution,» 11 12 2022. [Internett]. Available: [https://explore.maxar.com/Imagery-
\] Leadership-Spatial-Resolution](https://explore.maxar.com/Imagery-Leadership-Spatial-Resolution). [Funnet 17 02 2023].
- [27 Earthdata, *What is Remote Sensing?*, accessed 2023.
]
- [28 S.-E. Qian, «Enhancing space-based signal-to-noise ratios without redesigning the satellite,» SPIE, 05 01 2011.
] [Internett]. Available: [https://spie.org/news/3421-enhancing-space-based-signal-to-noise-ratios-without-
redesigning-the-satellite?SSO=1](https://spie.org/news/3421-enhancing-space-based-signal-to-noise-ratios-without-redesigning-the-satellite?SSO=1). [Funnet 04 03 2023].
- [29 Google Earth Engine, *API Reference - Earth Engine*, accessed on May 10, 2023.
]
- [30 e. a. Jed Matson, «Ultrasonic Flare Gas Flow Meter Techniques for Extremes of High and Low Velocity Measurement
] and Experience with High CO2 Concentration,» i *North Sea Flow measurement workshop*, Aberdeen , 2010.

- [31 European Space Agency , «Carbon dioxide monitoring satellite given the shakes,» 10 11 2021. [Internett]. Available:
] https://www.esa.int/Applications/Observing_the_Earth/Copernicus/Carbon_dioxide_monitoring_satellite_given_the_shakes. [Funnet 02 02 2023].
- [32 E. S. Agency, *Sentinel-5P Data Access and Products*, 2017.
]
- [33 Creodias, «Searching, processing and analysis of Sentinel-5P data on CREODIAS,» 02 01 2020. [Internett]. Available:
] <https://creodias.eu/searching-processing-and-analysis-of-sentinel-5p-data-on-creodias>. [Funnet 01 04 2023].
- [34 Copernicus, *Collection 0 Level-2A - Sentinel Online*, Accessed 2023.
]
- [35 C. Dempsey, «First Satellite Images from Sentinel-2 Delivered,» 29 06 2015. [Internett]. Available:
] <https://www.gislounge.com/first-satellite-images-from-sentinel-2-delivered/>. [Funnet 26 01 2023].
- [36 Satellite Imaging Corporation , «Landsat 8 Satellite Sensor (15m),» 01 02 2020. [Internett]. Available:
] <https://www.satimagingcorp.com/satellite-sensors/other-satellite-sensors/landsat-8/>. [Funnet 07 03 2023].
- [37 U. S. G. Survey, *Landsat 8*, 2023.
]
- [38 Maxar Technologies, 09 12 2022. [Internett]. Available: <https://blog.maxar.com/earth-intelligence/2022/mapping-methane-emissions-using-maxars-worldview-3-satellite>. [Funnet 03 02 2023].
- [39 European Space Agency, *WorldView-3*, 2023.
]
- [40 European Space Agency, *GHGSat*, 2023.
]
- [41 Earthdata, *NASA Selects GHGSat Data for Evaluation*, 2022.
]
- [42 «Aaron Clark and Janet Paskin,» Bloomberg, 06 11 2022. [Internett]. Available:
] <https://www.bloomberg.com/news/features/2022-11-06/-satellite-data-methane-release-climate-change?leadSource=uverify%20wal>. [Funnet 08 04 2023].
- [43 NASA, «Visible Infrared Imaging Radiometer Suite (VIIRS),» 02 03 2023. [Internett]. Available:
] <https://www.earthdata.nasa.gov/learn/find-data/near-real-time/viirs>. [Funnet 27 03 2023].
- [44 C. Cao, F. DeLuccia, X. Xiong, R. Wolfe og F. Weng, «Early on-orbit performance of the Visible Infrared Imaging Radiometer Suite (VIIRS) onboard the Suomi National Polar-orbiting Partnership (S-NPP) satellite,» *IEEE Trans. Geosci. Remote Sens*, 2013.

- [45 D. Wener, «Earth Science and Climate Monitoring | Researchers Turn to VIIRS for Tracking Oil Well Gas Flares,»
] Space News , 32 10 2013. [Internett]. Available: <https://spacenews.com/37788earth-science-and-climate-monitoring-researchers-turn-to-viirs-for-tracking/>. [Funnet 15 03 2023].
- [46 e. a. Alba Lorente, «Methane retrieved from TROPOMI: Improvement of the data product and validation of the first
] 2 years of measurements,» *Atmospheric Measurement Techniques* 14, pp. 665-684, 28 01 2021.
- [47 e. a. Evan D. Sherwin, «Single-blind validation of space-based point-source detection and quantification of onshore
] methane emission,» *Scientific Reports volume 13, Article number: 3836*, 07 03 2023.
- [48 e. a. ITZIAR IRAKULIS-LOITXATE, «Satellite-based survey of extreme methane emissions in the Permian basin,»
] *SCIENCE ADVANCES*, pp. Vol 7, Issue 27, 30 06 2021.
- [49 e. a. Can Lit, «Anewmachine-learning-based analysis for improving satellite-retrieved atmospheric composition
] data: OMISO2 asanexample,» *Atmos. Meas. Tech.*, 15, p. 5497–5514, 27 09 2022.
- [50 P. R. V. C. H. Y. W. A. G. M. W. F. H. C. M. P. B. G. R. W. C. a. B. H. Joyce, «Using a deep neural network to detect
] methane point sources and quantify emissions from PRISMA hyperspectral satellite images,» *EGUsphere [preprint]*,
<https://doi.org/10.5194/egusphere-2022-924>, 2022.
- [51 B. J. M. J. D. B. P. M. G. V. d. B. A.-W. P. S. L. A. B. T. H. S. V. D. J. M. J. J. D. G. M. I.-L. I. G. J. G. L. C. D. H. Schuit,
] «Automated detection and monitoring of methane super-emitters using satellite data,» *Atmos. Chem. Phys. Discuss.*
[preprint], 26 01 2023.
- [52 H. Hu, O. Hasekamp, A. Butz, A. Galli, J. Landgraf, J. Aan de Brugh, T. Borsdorff, R. Scheepmaker og I. Aben, «The
] operational methane retrieval algorithm for TROPOMI,» *Atmospheric Measurement Techniques*, vol. 9, p. 5423–
5440, 2016.
- [53 e. a. Otto Hasekamp, «Algorithm Theoretical Baseline Document for Sentinel-5 Precursor Methane Retrieval[1cm],»
] Netherlands Insitute for Space Research , 01 02 2019. [Internett]. Available:
<http://www.tropomi.eu/sites/default/files/files/publicSentinel-5P-TROPOMI-ATBD-Methane-retrieval.pdf>.
[Funnet 01 03 2023].
- [54 Flaring Monitoring Project , p. <https://www.flaringmonitor.org/index.html>, 15 04 2023.
]
- [55 Flaring Monitor, «Flaring Monitor User Interface,» Flaring Monitoring project , [Internett]. Available:
] <https://www.flaringmonitor.org/map/index.html?layer=3>. [Funnet 16 04 2023].
- [56 C. D. Elvidge, M. Zhizhin, F.-C. Hsu og K. E. Baugh, «VIIRS nightfire: Satellite pyrometry at night,» *Remote Sensing*,
] vol. 5, p. 4423–4449, 2013.
- [57 C. S. o. M. VIIRS Nightfire, *viirs-flare-data*, GitHub, 2021.
]
- [58 Google Earth Engine, *Sentinel-5P OFFL CH4: Offline Methane*, 2023.
]

- [59 C. Leyden, «Satellite data confirms Permian gas flaring is double what companies report,» *Environmental Defense Fund*, vol. 24, 2019.
- [60 E. D. Fund, *Permian Methane Analysis Project, PermianMAP*, 2021.
- [61 R. Schulz og T. Bredariol, *Flaring emissions – analysis*, IEA.
- [62 G. Plant, E. A. Kort, A. R. Brandt, Y. Chen, G. Fordice, A. M. Gorchov Negron, S. Schwietzke, M. Smith og D. Zavala-Araiza, «Inefficient and unlit natural gas flares both emit large quantities of methane,» *Science*, vol. 377, p. 1566–1571, 2022.
- [63 R. Kleinberg, «Greenhouse gas footprint of oilfield flares accounting for realistic flare gas composition and distribution of flare efficiencies,» *Authorea Preprints*, 2022.
- [64 J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai og others, «Recent advances in convolutional neural networks,» *Pattern recognition*, vol. 77, p. 354–377, 2018.
- [65 M. Etminan, G. Myhre, E. J. Highwood og K. P. Shine, «Radiative forcing of carbon dioxide, methane, and nitrous oxide: A significant revision of the methane radiative forcing,» *Geophysical Research Letters*, vol. 43, p. 12–614, 2016.
- [66 D. R. Caulton, P. B. Shepson, M. O. L. Cambaliza, D. McCabe, E. Baum og B. H. Stirm, «Methane destruction efficiency of natural gas flares associated with shale formation wells,» *Environmental science & technology*, vol. 48, p. 9548–9554, 2014.
- [67 *Methane Tracker 2020*, IEA, 2020.
- [68 *Global Methane Tracker 2022*, IEA, 2022.
- [69 J. Hu, Z. Yang og H. Su, «Dynamic Prediction of Natural Gas Calorific Value Based on Deep Learning,» *Energies*, vol. 16, p. 799, 2023.
- [70 Topcoder, «GRADIENT DESCENT IN MACHINE LEARNING,» 28 06 2019. [Internett]. Available: <https://www.topcoder.com/blog/gradient-descent-in-machine-learning/>.
- [71 MathWorks, *What Is Reinforcement Learning?*, accessed on April 14, 2023.
- [72 The Netherlands Space Office and AirbusDS , «TROPOMI TROPOspheric Monitoring Instrument,» 23 04 2023. [Internett]. Available: <http://www.tropomi.eu/data-products/methane>.
- [73 B. E. Mene, «Gas Flaring and Low Carbon Development: A Comparative Analysis of Nigeria, UK and Alberta,» *University of Calgary PRISM Repository*, 2019.

[74 World Bank, *Zero Routine Flaring by 2030 (ZRF) Initiative*, 2021.

]

[75 L. Sui, T. H. Nguyen, O. Khrakovsky, J. E. Matson, N. J. Mollo og M. P. Boespflug, «Ultrasonic Flowmeter for Accurately Measuring Flare Gas over a Wide Velocity Range,» 2009.

ATTACHMENTS

SENTINEL 5P METHANE SPECIFICATIONS

Data	Symbol	Units	Source	Pre-Process needs	If not available
S5P level 1b Earth radiance SWIR band	I	mol/s/m ² / nm/sr	S5P level 0-1b product	per ground pixel	no retrieval
S5P level 1b Earth radiance NIR band	I	mol/s/m ² / nm/sr	S5P level 0-1b product	per ground pixel, spatially co-located with SWIR ground pixel	no retrieval
S5P level 1b Solar irradiance SWIR band	$F_{0,meas}$	mol/s/m ² / nm	S5P level 0-1b product		use previous measurement
S5P level 1b solar irradiance NIR band	$F_{0,meas}$	mol/s/m ² / nm	S5P level 0-1b product		use previous measurement
latitude	lat	degree	S5P level 0-1b product		no retrieval
longitude	lon	degree	S5P level 0-1b product		no retrieval
solar zenith angle	θ_0	degree	S5P level 0-1b product		no retrieval
viewing zenith angle - SWIR band	θ_v	degree	S5P level 0-1b product		no retrieval
relative azimuth angle –SWIR band	ϕ	degree	S5P level 0-1b product		no retrieval
viewing zenith angle - NIR band	θ_v	degree	S5P level 0-1b product		no retrieval
relative azimuth angle - NIR band	ϕ	degree	S5P level 0-1b product		no retrieval
cloud fraction for S5P SWIR and NIR ground pixel			VIIRS / RAL algorithm	per ground pixel, spatially co-located with SWIR and NIR ground pixel.	use FRESCO apparent pressure and backup cloud screening
cirrus reflectance for S5P SWIR and NIR ground pixel		[-]	VIIRS / RAL algorithm	per ground pixel, spatially co-located with SWIR and NIR ground pixel.	use FRESCO apparent pressure and backup cloud screening
apparent pressure for ground pixel and surrounding		Pa	FRESCO (L2)	find corresponding and neighbouring ground pixel	filter based on retrieved scattering parameters
non scattering retrieval results for weak and strong CH ₄ and water bands for ground pixel and surroundings	$[CH_4]_{weak}$ $[CH_4]_{strong}$ $[H_2O]_{weak}$ $[H_2O]_{strong}$	mol cm ⁻²	CO algorithm SICOR (L2)	Find corresponding and neighbouring ground pixel	filter based on retrieved scattering parameters

Continued on next page

FLARING MONITOR – USER INTERFACE

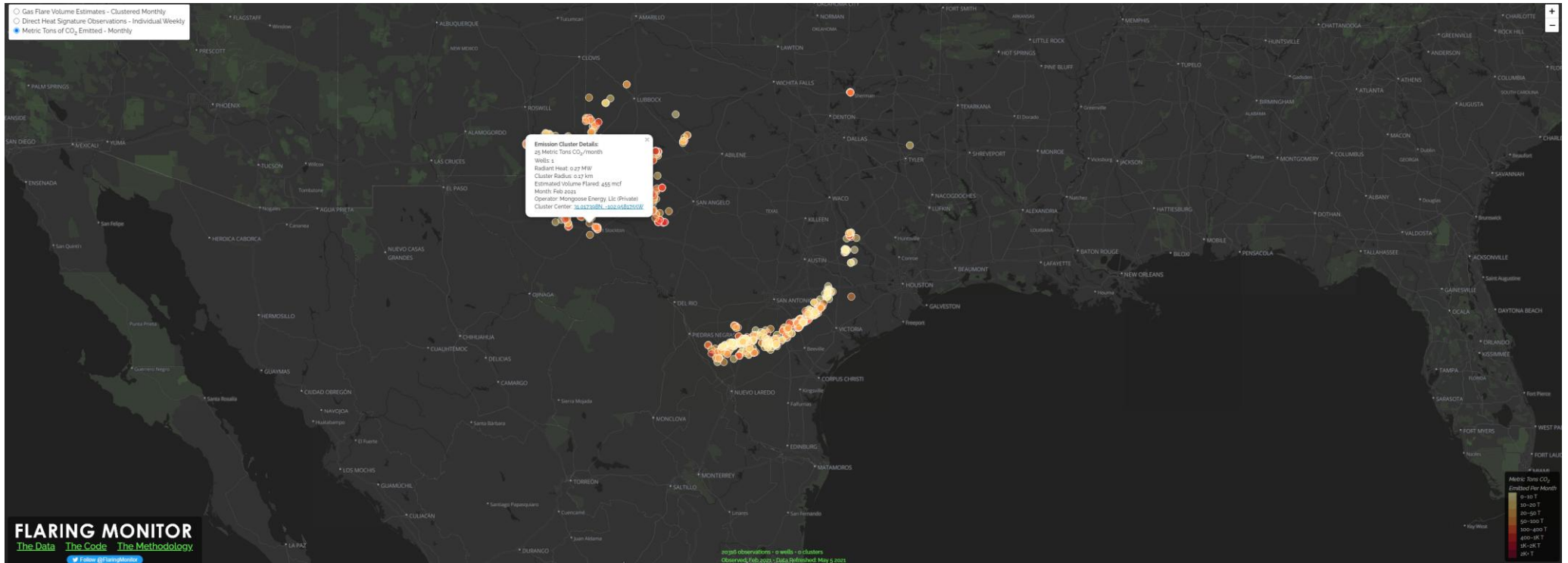


Figure 42 Detailed overview of the Flaring Monitoring user interface that provides CO2 emission data from US oil and gas sites based VIIRS data



Norges miljø- og biovitenskapelige universitet
Noregs miljø- og biovitenskapelige universitet
Norwegian University of Life Sciences

Postboks 5003
NO-1432 Ås
Norway