

Article

Not peer-reviewed version

Investigating the Performance of Open Vocabulary Classification Algorithms for Pathway and Surface Material Detection in Urban Environments

[Kauê de Moraes Vestena](#)*, [Silvana Phillipi Camboim](#), [Maria Antonia Brovelli](#), [Daniel Rodrigues dos Santos](#)

Posted Date: 17 September 2024

doi: 10.20944/preprints202409.1321.v1

Keywords: Open-Vocabulary Algorithms; Pavement Segmentation; Surface Material Detection; Street-View Imagery; Urban Mobility



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Article

Investigating the Performance of Open Vocabulary Classification Algorithms for Pathway and Surface Material Detection in Urban Environments

Kauê de Moraes Vestena ^{1,*}, Silvana Phillipi Camboim ¹, Maria Antonia Brovelli ²
and Daniel Rodrigues dos Santos ²

¹ Federal University of Paraná

² Politecnico di Milano

³ Military Institute of Engineering

* Correspondence: kauemv2@gmail.com

Abstract: Mapping pavement types, especially in sidewalks, is essential for urban planning and mobility studies. Identifying pavement materials is a key factor in assessing mobility, such as walkability and wheelchair usability. However, satellite imagery in this scenario is limited, and in-situ mapping can be costly. A promising solution is to extract such geospatial features from street-level imagery. This study explores using open vocabulary classification algorithms to segment and identify pavement types and surface materials in this scenario. Our approach uses Large Language Models (LLMs) to improve the accuracy of classifying different pavement types. The methodology involves two experiments: the first uses free prompting with random street-view images, employing Grounding Dino and SAM algorithms to assess performance across categories. The second experiment evaluates standardized pavement classification using the Deep-Pavements Dataset and a fine-tuned CLIP algorithm optimized for detecting OSM-compliant pavement categories. The study presents open resources, such as the Deep Pavements Dataset and a fine-tuned CLIP-based model, demonstrating a significant improvement in the True Positive Rate (TPR) from 56.04% to 93.5%. Our findings highlight both the potential and limitations of current open vocabulary algorithms and emphasize the importance of diverse training datasets. This study advances urban feature mapping by offering a more intuitive and accurate approach to geospatial data extraction, enhancing urban accessibility and mobility mapping.

Keywords: open-vocabulary algorithms; pavement segmentation; surface material detection; street-view imagery; urban mobility

1. Introduction and Related Work

The study of sidewalks is crucial for accessibility and quality of life in cities [1–4]. As the required level of detail is not available in satellite imagery and in-situ surveys are costly, mapping such features is challenging. Using street-level imagery for this purpose is, therefore, a potential improvement. However, defining the features to be classified in street-level imagery is challenging because landscape features are conceptualised in natural language. Thus, Large Language Models (LLMs) can be used to improve the process of image pattern retrieval from street-level imagery by allowing the use of flexible and natural language prompts. This characteristic makes the extraction of map features more intuitive and accurate by exploiting the language-understanding capabilities of LLMs, which can be applied in many contexts, such as in urban planning and transport studies, such as in [5], where it was used to solve a community-level land-use task through using a feedback iteration. Recently, there has been considerable interest in the accessibility of these models, as

demonstrated by ChatGPT, which serves as an interface for user interaction with the latest iterations of the Generative Pre-trained Transformer (GPT). This development is the most recent in a continuing series of advances in natural language processing. The study of Natural Language Processing (NLP) can be traced back to the 1950s [6], with early developments culminating in the establishment of the essential subtasks, such as sentence boundary detection, tokenisation, part-of-speech assignment to individual words, morphological decomposition; chunking; and problem-specific segmentation [6]. Later developments resulted in the tailoring of conventional algorithms, such as Support Vector Machines and Hidden Markov Models.

Among these developments, the extension of open vocabulary prompting for image pattern retrieval represents a substantial leap forward in NLP, driven by advances in LLM. This approach has recently gained more attention due to its capabilities of richer expressiveness and more extensive flexibility due to generalisation capabilities [7], its ability to capture nuances of processes [8], and its enabling for transfer-learning [9]. The work of Zareian et al. [10] claims to be the first to pose the problem of "Open Vocabulary Object Detection" in opposition to other similar approaches. One example of that alternative approach was the "Zero Shot-Detection," which aims to generalise from seen to unseen classes using semantic descriptions or attributes [11–13]. Another would be Weakly Supervised Detection, which focuses on detecting classes with limited, such as image-level or noisy labels [14,15]. It is worth noticing that there are some main drawbacks to Open-vocabulary approaches, such as a higher computational cost due to their larger complexity [7] and the risk of misinterpretation due to the lack of constraints [8].

This paper investigates the use of open vocabulary algorithms for the domain-specific task of segmenting pavement regions along their surface materials. Identification of pavement material types is crucial for maintaining road safety and ensuring the well-being of people [16]. This process is vital for mobility studies as it influences safety, skid resistance, and road noise [17]. Particularly for sidewalks, identifying suitable pavement types is essential for improving accessibility, urban mobility, and safety for all users, including those with reduced mobility [18]. Additionally, using specific pavement types like exposed aggregate concrete (EAC) and porous concrete can significantly reduce noise and enhance safety on both roadways and sidewalks [19]. Zeng & Boehm [20] exploited a broader investigation of open-vocabulary algorithms, which got averagely good results compared to SOTA closed-vocabulary ones, observing a prompt-dependent accuracy.

Although "street" and "sidewalk" are common classes amongst many classification/segmentation datasets [21–24], being one of the reference classes for testing many proposed algorithms [24–26], pavement-type identification is a far scarcely undertaken task: Convolutional Neural Networks were used to identify few classes limited to asphalt, gravel, and cement [27], and also for the identification of asphalt damage traits such as "pothole" or "patch" [28], and pavement damage assessment [29], pixel-level segmentation, limited to "paved" and "unpaved" streets [30], and only recently including eight specific classes such as "granite blocks" and "hexagonal asphalt paver" [31]. Each of these studies developed specific models with a closed set of classes, thus not allowing a wide use of the algorithms, in contrast to open vocabulary counterparts such as CLIP [32] and Grounding Dino [33], designed to be generalist algorithms.

Furthermore, none of these previous studies on pavement material identification follows a global standard. This work innovates by integrating its application into OpenStreetMap (OSM), the world's most popular open and collaborative geographic database[34]. The feature attribute standards created by the OSM community have typically been agreed upon over many years, including surface materials, although they are often not defined by domain experts. Adopting this strategy helps maintain interoperability between the developments of this research and other tools and applications linked to OSM. In this context, this work aims to test the behaviour of some state-of-the-art algorithms for open vocabulary image classification tasks in identifying pavement types according to the OSM conceptual model, thus aiming to help design valuable building blocks for producing compliant geographic data.

This study presents a valuable proposition to the community of urban planners and geospatial analysts. By using open vocabulary algorithms and large language models (LLMs) to segment and

identify pavement materials in street-level imagery, this work advances the methodology for mapping urban features. It bridges the gap between natural language processing and geospatial data extraction, providing a more intuitive and granular approach to feature classification. By integrating with OpenStreetMap (OSM), the study improves data interoperability and sets a precedent for using open and collaborative platforms for urban data management. These innovations contribute to more accessible, scalable and flexible solutions for urban infrastructure analysis, with potential applications in transport planning, accessibility improvement and public safety.

2. Methodology

The task of using NLP to extract features from images has as its main challenge the process called "grounding", whose level means the capability of the model to associate elements of language with the proper regions in the images [35–37]. It was shown that the relevance of the hierarchical relationships between concepts improves the grounding level by establishing ontologies [38–40]. Ontologies bridge the gap between natural language and logical reasoning by providing a machine-readable representation (class) of real-world concepts [41]. Our study's methodological approach, presented at Figure 1 with its main steps and processes, aimed at fulfilling our study's goals.



Figure 1. study's methodology workflow.

The motivational task behind this study can be formalised as follows: "Given a set of street view images, the objective is to segment all visible paths in each image using free-text input algorithms, identify their surface materials, and ensure that the process is not affected by the use of synonyms in the input." Considering the hierarchical relationships between the entities of interest and their

properties, an ontology was created using the Web Ontology Language (OWL) standard [42]. This ontology is presented as a graph in Figure 2. The hierarchy naturally starts with "pathway," which is then specialised into "Road" and "Footway." The branching follows up to two main interest categories: "Road" and "Sidewalk." All of them, at the same time, have their main differences. Still, they fundamentally share the characteristic of being walkable, which is essentially made possible through having a surface material, hereby treated as a fundamental shared characteristic of all "pathways". After that, "surface material" has its branching sub-ontology that is composed of two main categories: "unpaved material" and "paved material," the latter being subdivided into "continuous" and "block-based". With the ontology set, it is then possible to establish different degrees of semantic separation among classes, highlighting pavement surfaces. Therefore, "asphalt" is more similar to "concrete" (0 degrees, same hierarchical level) than it is to "sett" (1 degree) than it is to compacted pavement (2 degrees). This ontology is basilar to the proposed experiment, with its hierarchical relationships being called of along the study. As far as our knowledge goes, no similar ontology was proposed.

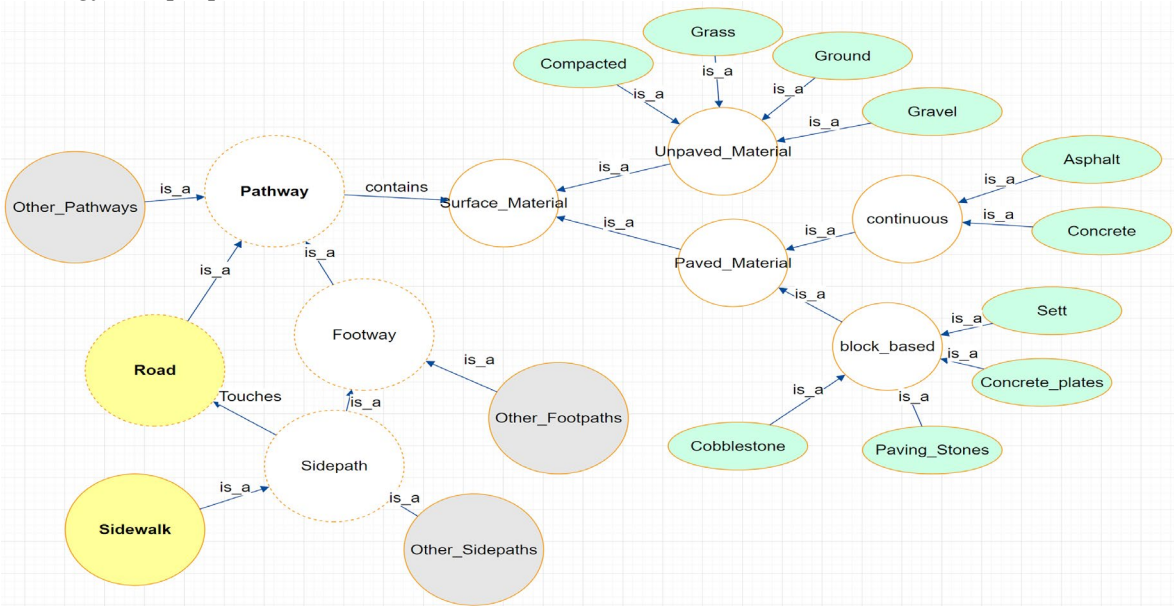


Figure 2. the ontology of the interest phenomena.

To undertake the present study, we carried out two different experiments. In the first experiment, named "Free Prompting and Random Imagery Evaluation," we tested free prompting against random imagery to observe the actual outcomes when the prompts are used openly, akin to real-usage scenarios aiming for mass processing. The employed algorithm was a combination of Grounding Dino [33] and the foundation model SAM - Segment Anything [43]. Grounding Dino provided output bounding boxes based on the prompts, which were then used as inputs for SAM to obtain the segmented features. We used Mappillary, an ever-growing platform with billions of community-contributed CC BY-SA licensed street-view images [44] for the data source.

In the second experiment, "Evaluation of Open-Vocabulary Algorithms for Standardized Pavement Classification," we tested standardised pavement classes prompting against images of surface patches to assess the classification ability of an open-vocabulary algorithm in a constrained setup. We employed the CLIP algorithm, known for assigning probabilities to each entry of a set of sentences, that can be changed at each query, constituting a fundamental advantage for it allows adaptation to different realities. We tested eight variations of the algorithm, each employing different datasets and training strategies, as compared in Table 1 regarding their origin and computational costs. The data source for this experiment was the "deep-pavements-dataset," a tailored dataset of surface patches collaboratively maintained on GitHub [45], containing 500 samples for OSM-

compliant pavement categories at the time of the experiments. The testing included all categories in the ontology, with a few samples shown in Figure 3. Future releases may include more categories.

Table 1. tested CLIP algorithm variants.

Package	Model Name	Base Algorithm	Main Image Dataset	Model Size (GB)	Image	Text	Image GFLOPS	Text GFLOPS	Package	Model Name
					Parameters (Millions)	Parameters (Millions)				
CLIP	RN101	CNN	OpenAI	0.28	56.26	63.43	19.54	5.96	CLIP	RN101
CLIP	RN50x64	CNN	OpenAI	1.26	420.38	202.88	529.11	23.55	CLIP	RN50x64
CLIP	ViT-B/32	VT	OpenAI	0.34	87.85	63.43	8.82	5.96	CLIP	ViT-B/32
CLIP	ViT-L/14	VT	OpenAI	0.89	303.97	123.65	162.03	13.3	CLIP	ViT-L/14
Open CLIP	ViT-H-14-378-quickgelu	VT	dfn5b	3.95	632.68	354.03	1006.96	47.09	Open CLIP	ViT-H-14-378-quickgelu
Open CLIP	EVA02-E-14-plus	VT	laion2b_s9b_b1 44k	10.1	4350.56	694.33	2264.33	97.86	Open CLIP	EVA02-E-14-plus



Figure 3. Some snapshots of the deep-pavements-dataset.

It is valuable to elaborate on the reasons for selecting this set of algorithms for our study. SAM is particularly noted for its unpaired segmentation capabilities, supported by a billion-level sample size, which enhances its applicability to real-world scenarios and provides superior generalisation performance [46,47]. Grounding Dino offers complementary strengths and significantly outperforms comparable algorithms. This advantage is attributed to its use of a transformer architecture that integrates multi-level text information [33,48]. CLIP marked a significant breakthrough in the vision-language domain by employing a shared embedding space for text and images created through a contrastive learning approach [49]. It has been recognised for its robustness across various scenarios [50–52]. These algorithms are considered zero-shot learners capable of performing tasks not specifically optimised [11]. Furthermore, despite their relatively recent introduction, these models have already seen widespread use in the industry for applications such as automated image data annotation, image search engines, accessibility tools providing image descriptions, and enhancing content recommendation systems [53].

Regarding the first experiment and the proposed ontology, we tested the following categories with their corresponding prompts:

1. Auxiliary: These entities are detected to determine if a pavement detection failure occurred due to occlusion rather than incorrect detection. Prompts are "car", "vehicle", "pole", "tree".

2. Sidepaths: This query is primarily aimed at detecting sidewalks, but it may also retrieve other sidepaths, such as paved shoulders. Prompts are "sidewalk", "sidepath", "sideway", "sidetrack", and "lateral".
3. Roads: Focused on detecting motorised pathways. Prompts are "road" and "street".
4. Pathways: These are intended to detect any kind of traversable way. Prompts are "way", "path", "pathway", "pavement", and "track".
5. Surface Pavement Types: Prompts directly target the property, with additional ones included for broader testing. Prompts are "sett", "grass", "cobblestone", "earth", "soil", "dirt", "sand", "concrete", "paving stones", "chip seal", "gravel", "compacted", "asphalt", "concrete plates", and "ground".
6. Abstract Concepts: Words that do not have a unique or physical representation are used to test the model's responses. Prompts include anything, nothing, something, void, and thing.

For the second experiment, we conducted two tests with different strategies to compare the performance with the pre-trained models. In the first one, called the "Model Combining" test, we summed the class probabilities across all models to create a score, selecting the class with the highest score as the winner, thus just using. We performed this test using all models and with the three best overall scorers.

Second, in the so-called "Fine-Tuning" test, we selected 60% of the samples for fine-tuning and used the remaining 200 samples per category for testing. Following the specifications in Table 1, we chose the lightest model to evaluate the extent of improvement relative to the model size and computational burden. Previous studies [49] show that CLIP's fine-tuning is highly sensitive to the optimiser's algorithm hyperparameters. Therefore, we empirically tested multiple scenarios, presenting the worst and the best outcomes side by side. All the produced models and analyses are published in a repository at the HuggingFace community [54] to ensure reproducibility.

3. Results and Discussion

3.1. Free Prompting and Random Imagery Evaluation

Following the methodology presented in the previous section, eight images were tested for each prompt. So, for categories that received five prompts, the experiment involved 40 images. Table 2 presents the results obtained for the "auxiliary" category. In the forthcoming tables, a "***" sign indicates that the targeted prompt does not exist in the image sample; thus, the correct outcome should be none.

Table 2. prompt results for the "auxiliary" category.

row	Prompt	1	2	3	4	5	6	7	8	hit %
I	car	car	truck	car	car	car panel*	car	car	car	100
II	poles	street lamp	power pole	building pole	street lamp	street lamp	guard rail **	street lamp	street lamp	87.5
III	vehicle	car	truck	car	truck	fence**	car	car	truck	87.5
IV	tree	tree	tree	tree	tree	tree	vegetatio n	tree	tree	100

The desired results were achieved in all cases, except when the targeted feature was absent in the picture, leading the algorithm to hallucinate. The Table 3 presents the results for the "Sidepaths" category.

Table 3. prompt results for "sidepaths" category.

row	Prompt	1	2	3	4	5	6	7	8	hit %
-----	--------	---	---	---	---	---	---	---	---	-------

I	sidewalk	sidewalk	sidewalk	road	sidewalk	sidewalk	a sidepath	kerbs	sidewalk	75
II	sidepath	road **	road **	kerb	road **	road	road **	road	road	0
III	sideway	car panel	car	truck	car	car	car panel	road	car	0
IV	sidetrack	car	car panel	building	car	car panel **	truck container	car	car	0
V	lateral	road **	road	road and sidewalk	car **	car	vegetatio n all pavemen ts	road	car **	0
VI	walkway	road and sidepath	road **	road **	road **	partial sidewalk		road	road **	12.5

There is considerable evidence supporting the hypothesis of a lack of generalisation capability. While "sidewalk" produced consistently good results, none of the other prompts were recognised as similar concepts, which should have led to proper detections. It's worth recalling that "sidewalk" is a common category in many well-spread training datasets. Table 4 presents the results obtained for the "roads" category, reinforcing the hypothesis of a stronger capability for road identification, likely due to the prevalence of roads in terrestrial imagery.

Table 4. prompt results for the "sidepaths" category.

row	Prompt	1	2	3	4	5	6	7	8	hit %
I	street	road	road	road	road	road	road	road	road	100
II	road	road	road	road	road	road	road	road	road	100

Table 5 presents the results for the "pathway" category. Although the results were generally good, there was a noticeable bias towards detecting roads.

Table 5. prompt results for the "pathway" category.

row	Prompt	1	2	3	4	5	6	7	8	hit %
I	track	car	road	road	road	road	road	road	road	87.5
II	pavemen t	road	road	a path	road	road --	road	walkway	road	100
III	way	bus	road markings	road	road	bus	road	sidewalk	road	62.5
IV	path	road	road	road	road	road	road	road	road	100
V	pathway	road	road	sidewalk	road	road	road	road	road	100

Although, on average, we can see the expected results, there's a clear bias towards the detection of roads. Table 6 presents the results for the "surface" category, where only "asphalt" and "grass" yielded satisfactory results. Other prompts, such as "sand" and "concrete," were less successful due to the nature of the test, which involved random images that often did not correspond to the prompts. The correct behaviour of returning "nothing" was only observed in the "concrete plates" category.

Table 6. prompt results for the "surface" category.

row	Prompt	1	2	3	4	5	6	7	8	hit %
I	sett	car **	car **	car **	car **	asphalt road **	tree **	car **	car **	0
II	grass	grass sidepath	grass sidewalk	grass sidewalk	grass in a garden	grass sidepath sidepath	plants **	grass stripe	grass sidepath	87.5
III	cobblesto ne	grass sidepath	a pathway with stones	bush row	dirt sidepath	with vegetatio n	asphalt road	asphalt road	asphalt road	0
IV	earth	car **	gravel sidewalk **	asphalt road **	car panel **	asphalt road **	a tire **	asphalt road **	car **	0

V	soil	an unpaved sidepath	asphalt road **	asphalt road	soy crop	car	soil sidepath	grass sidewalk	grass sidepath	12.5
VI	dirt	asphalt road **	asphalt road **	car **	unpaved sidewalk **	asphalt road **	asphalt road and car panel **	a grass slope **	asphalt road **	0
VII	sand	car **	asphalt road **	asphalt road **	sand sidewalk	asphalt road **	asphalt road **	asphalt road (partial) **	car	12.5
VIII	concrete	asphalt road **	compacted pathway **	asphalt road **	asphalt road **	concrete sidewalk	guard rail **	asphalt road **	asphalt road **	12.5
IX	paving- stones	asphalt road	soil sidepath **	concrete sidewalk **	asphalt road	asphalt road **	asphalt road	asphalt cobblestone sidewalk stripe	asphalt road	0
X	chipseal	car	truck	car	car	car	car	car	car	0
XI	gravel	asphalt road **	asphalt pathway **	asphalt pathway **	sidepath with vegetatio n	car	asphalt road **	car	wall	0
XII	compacte d	truck	car	car	car	car	car	car	car	0
XIII	asphalt	asphalt road	asphalt road	asphalt road	asphalt road	asphalt pathway	asphalt road	asphalt road	asphalt road	100
XIV	Concrete plates	truck	nothing	nothing	truck	truck	truck	truck	a stairway	25
XV	ground	asphalt road **	asphalt road **	asphalt road and ground sidewalk	asphalt road and asphalt sidepath	asphalt road **	asphalt road **	asphalt road	asphalt road **	0

We got suitable matches only for 'asphalt' and 'grass'. Some, such as "sand" and "concrete," were affected by the nature of the test, where random images were selected and often had no correspondence. Although the correct behaviour should be "nothing," this was only seen in the "concrete plates" category. Table 7 presents the "abstract" category results, where a bias towards vehicles was observed. The prompt "anything" yielded more varied answers than "something," while "void" and "nothing" still produced responses, indicating that the language model struggles with abstract concepts.

Table 7. prompt results for the "abstract" category.

row	Prompt	1	2	3	4	5	6	7	8	hit %
I	anything	car	landscape	road	person	road and sky	car panel	vegetation and sky	road and car panel	N/A
II	nothing	asphalt road and noise	asphalt road and noise	car	asphalt road and noise	noisy asphalt road and sky	sky and noise	truck	road, wall, and noise	N/A
III	something	car hood	car panel	traffic signal	car	car	nothing	truck	truck	N/A
IV	anything	road, sidewalk, and vegetation	noisy asphalt road and noisy sky	car	tree	landscape and noisy road	nothing	car	building	N/A

V	void	asphalt road	car	car	nothing	car	car	car	car	N/A
---	------	--------------	-----	-----	---------	-----	-----	-----	-----	-----

Once again, the bias towards vehicles can be observed. Somehow, "anything" gives more varied answers than "something". "Void" and "nothing" still produce matches, so the language model is not handling this abstract concept correctlyIn Figure 4, which summarises the results of the first experiment, we have applied the same concepts of generalisation through hierarchical categorisation as in the ontology in Figure 3.

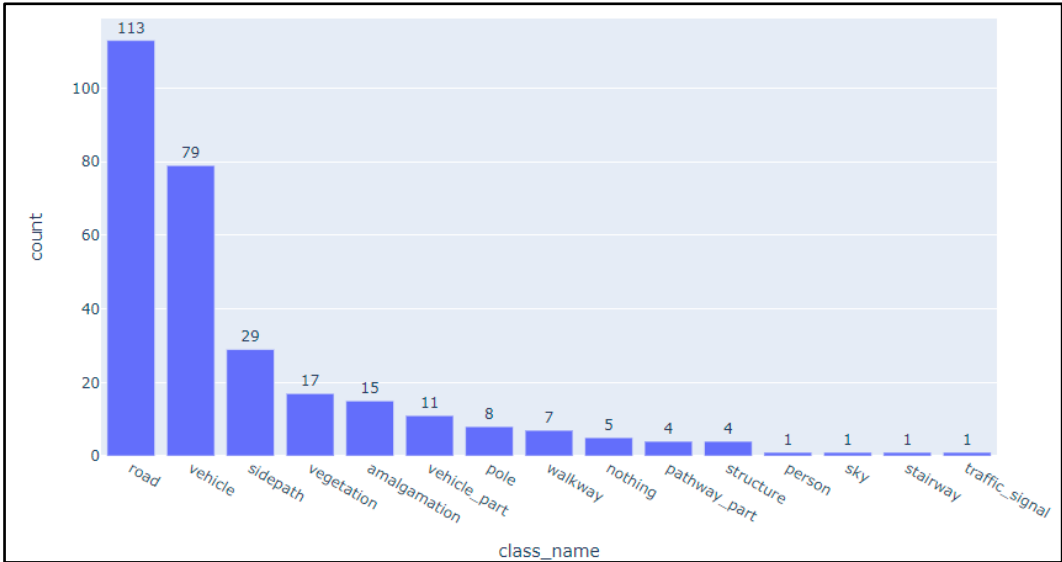


Figure 4. test 1 summarized results.

Even though "sidepath" is the third most common category, it is observable a bias towards roads and vehicles.

3.2. Evaluation of Open-Vocabulary Algorithms for Standardized Pavement Classification

We mainly relied on confusion matrices for the second experiment to summarize the algorithms' performances. In Figure 5, we computed those confusion matrices for all the presented versions.

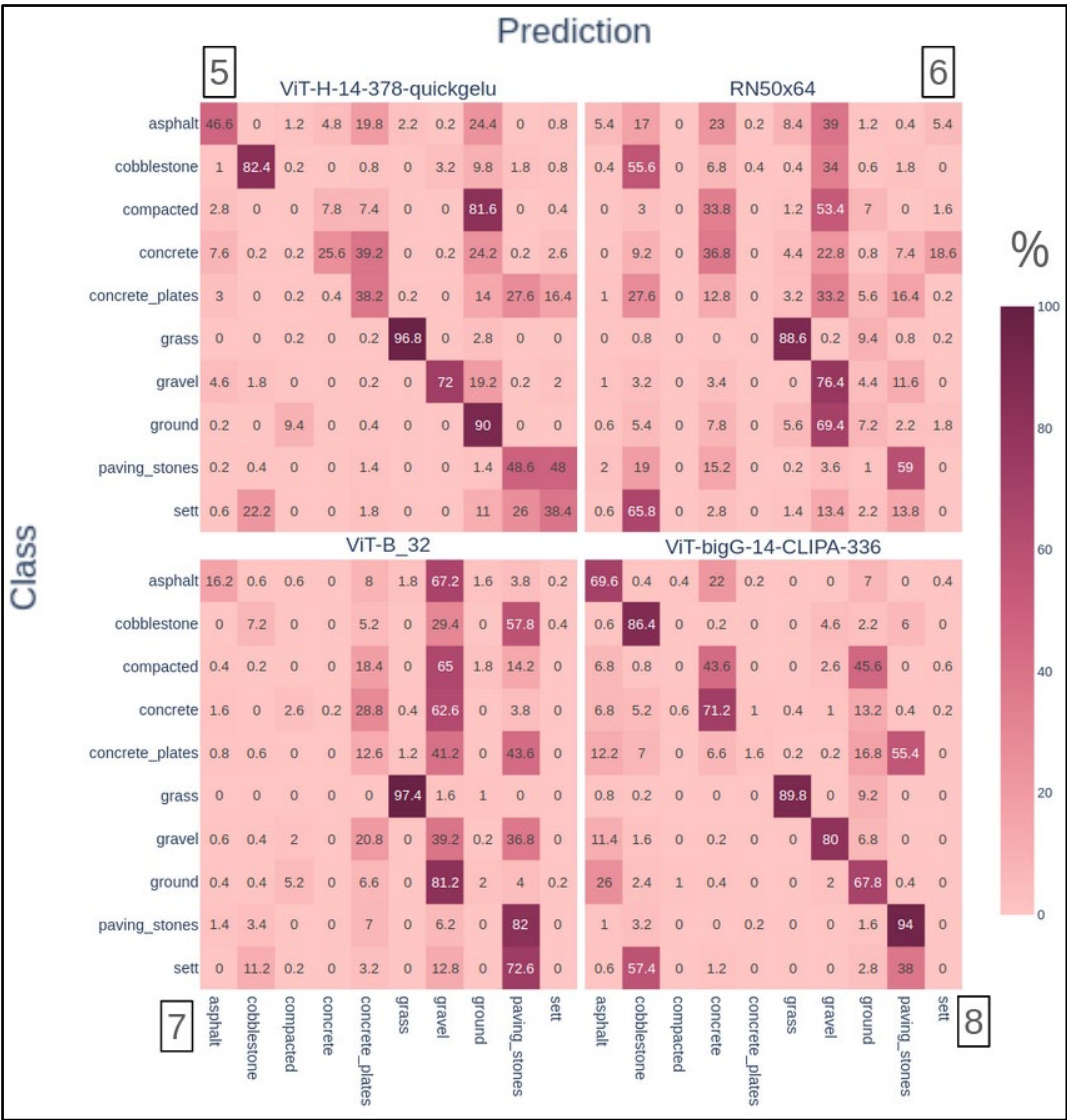


Figure 5. the resulting confusion matrices for the variants of the algorithm.

All True Positives are located on the main diagonal in the confusion matrices, with the rows summing to 100%. The columns reveal biases towards certain classes, such as the "Sett" category, which had notably low prediction rates (e.g., only 38% True Positives in model 5 and 0% in model 3). Extreme cases were observed in model 7, where the "gravel" class summed to 406% and the "ground" class to 318%, indicating significant misclassification issues. It's worth recalling that different datasets were used in each model's training. Therefore, both extremes can be due to imbalanced datasets, In this context, despite their differences, all models are variations of an algorithm that claims to be widely capable of generalising correctly both in vision and language domains. It not only far from happening but also occurs in very different ways throughout all models.

3.3. Combine and Fine-Tune Models

The results of the combining strategy, shown in Figure 6, that were obtained by summing the class probabilities from all models to create a composite score, with the highest-scoring class being selected as the winner. This test was conducted using all models and the three best overall performers.



Figure 6. confusion matrices for the combination strategy.

While a general improvement can be noticed, some classes still suffer from mixing with others: "sett" with "cobblestone" and "paving stones", "concrete plates" with "paving stones" and "ground", and compacted with concrete and ground. Some classes are still completely hindered, mainly "sett", "concrete plates" and "compacted".

Considering the fine-tuning experiments in the Figure 7, the results of the fine-tuning strategy are presented, and in the Table 8 the different settings for the optimizer hyperparameters are presented. The finetuning and, as a consequence, the selection of the best model were both made empirically, as reportedly done due to the complexity of AI models and the general lack of a closed-form solutions [55]. As stated in the methodology, we selected the model ViT-B/32 (number 7) as it got the worst overall performance, as later shown in Figure 8.

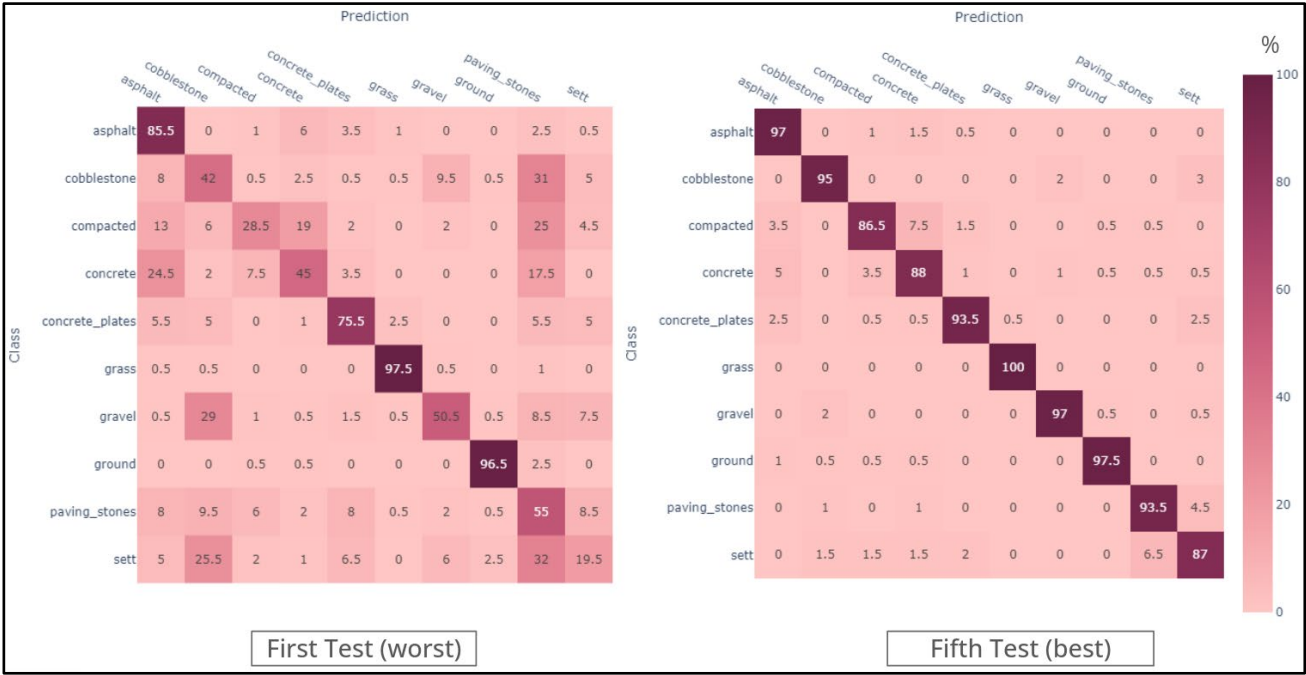


Figure 7. confusion matrices for the fine-tuning strategy.

Table 8. hyperparameter variation on finetuning results.

Hyperparameter	Original	Best Result
Batch size	256	320
Learning Rate	5.00E-05	5.00E-06
Weight Decay	0.2	0.5
AMSGrad Method	Deactivated	Activated
Betas	(0.9,0.98)	(0.9,0.98)
Epsilon	1.00E-06	1.00E-06
training epochs	100	100

The Figure 7 Shows that the improvement was quite prominent after reducing the learning rate by a whole order of magnitude, which allowed for better avoidance of local minima, as foretold in the literature [56], the use of the AMS grad [57] variant also gave adaptive capabilities, enabling the variation of the learning rate when mathematically viable. It performs slightly better than the state-of-the-art performance, which, as of April 2024, is an average TPR—True Positive Rate of 92.4% on ImageNet [58], according to [59]. It's worth recalling that the best improvement comparison of the final result is with image 5, confusion matrix number 7, which was the starting spot. To compare all the tested scenarios directly, we created the chart in Figure 8 Figure, containing the average TPR and the standard deviation for each tested variation.

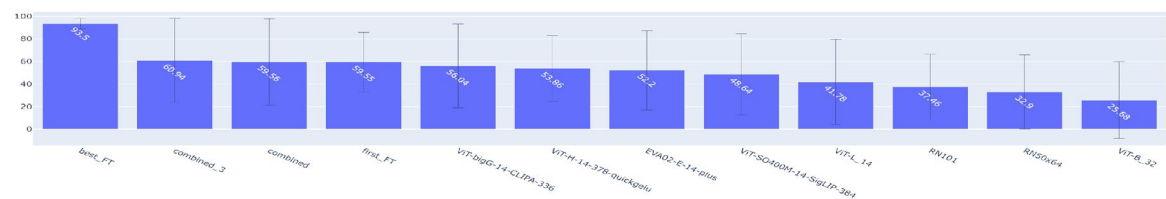


Figure 8. Average TPR of the different variants of the algorithm.

The chart in Figure 8 was created to directly compare all tested scenarios, containing the average TPR and the standard deviation for each tested variation. Enforcing a specific meaning for a word with many meanings, such as "compacted," played a central role in achieving these groundbreaking results, which are also very consistent among classes. Additionally, defining "sett" as quarried and roughly parallelepipedal bricks was essential, as image searches for "cobblestone" often yielded results that matched this description, reflecting an algorithmic bias corroborated by Figure 5, and Figure 6.

In summary, the first experiment's "Free Prompting and Random Imagery Evaluation" demonstrated that while the desired results were generally achieved, the algorithm struggled with hallucinations when targeted features were absent. This result was particularly evident in the "sidepaths" and "surface" categories, where a bias towards detecting roads was observed. The second experiment, "Evaluation of Open-Vocabulary Algorithms for Standardized Pavement Classification," highlighted the strengths and weaknesses of different CLIP algorithm variants. Confusion matrices revealed biases towards certain classes and the challenges in accurately classifying less common pavement types. Fine-tuning and combining strategies showed improvements, but some classes remained problematic. These findings underscore the importance of diverse training datasets and the potential for further enhancements in open-vocabulary models.

4. Conclusions

This study makes key contributions to geospatial analysis and urban infrastructure management. Firstly, we created the Deep Pavements Dataset, a robust and expandable dataset specifically designed for training models in standardised (OSM-compatible) pavement material

classification. This open dataset is a valuable resource for future research, enabling more comprehensive training and benchmarking of algorithms in this specialised area.

Secondly, we developed a fine-tuned, lightweight CLIP-based model optimised for pavement type detection. This model demonstrates how adapting state-of-the-art large language models to specific, real-world applications can lead to improved accuracy and efficiency. It also highlights the potential of using LLMs for tasks that require nuanced understanding and classification, offering a more intuitive approach to complex geospatial data challenges.

Our findings emphasise the importance of diversifying training datasets to improve the performance of open-vocabulary models in specialised tasks like pavement type detection. This result aligns with broader challenges in AI, where the reliance on widely used datasets may lead to skewed or limited results in less common applications. Addressing these biases by creating more diverse and representative datasets can significantly enhance model performance and reliability.

Future research should focus on enhancing the CLIP-based model by incorporating different or additional approaches, such as data augmentation and fine-tuning with heavier versions, to compare their performance on generic tasks. Exploring robust solutions that combine specialised CLIP models with constrained vocabularies to verify hallucinations independently could further enhance the reliability of these models. Achieving a truly generalist open-vocabulary algorithm for specialised tasks like pavement type detection is an ambitious goal. Still, it is essential for advancing the capabilities of AI in diverse and practical applications.

In conclusion, this study contributes to developing specialised AI models for urban planning and geospatial analysis and provides valuable insights into the broader implications of dataset biases and model generalisation. These contributions lay the groundwork for future advancements in creating more flexible, accurate, and representative AI models for real-world applications.

Supplementary Materials: The deep-pavements dataset is widely available at https://github.com/kauevestena/deep_pavements_dataset; the finetuned CLIP-based model is widely available at <https://huggingface.co/kauevestena/clip-vit-base-patch32-finetuned-surface-materials>.

Author Contributions: Conceptualization, Kauê de Moraes Vestena, Silvana Phillipi Camboim, Maria Brovelli and Daniel Rodrigues dos Santos; Data curation, Kauê de Moraes Vestena; Formal analysis, Kauê de Moraes Vestena; Funding acquisition, Silvana Phillipi Camboim and Maria Brovelli; Investigation, Silvana Phillipi Camboim; Methodology, Kauê de Moraes Vestena, Silvana Phillipi Camboim, Maria Brovelli and Daniel Rodrigues dos Santos; Project administration, Silvana Phillipi Camboim; Resources, Silvana Phillipi Camboim and Maria Brovelli; Software, Kauê de Moraes Vestena; Supervision, Silvana Phillipi Camboim, Maria Brovelli and Daniel Rodrigues dos Santos; Validation, Kauê de Moraes Vestena; Visualization, Kauê de Moraes Vestena; Writing – original draft, Kauê de Moraes Vestena; Writing – review & editing, Silvana Phillipi Camboim, Maria Brovelli and Daniel Rodrigues dos Santos. All authors have read and agreed to the published version of the manuscript.

Funding: Please add: This study was financed in part by the Coordination for the Improvement of Higher Education Personnel – Brazil (CAPES) – Finance Code 001.

Data Availability Statement: deep-pavements dataset is available under a MIT license: https://github.com/kauevestena/deep_pavements_dataset/blob/main/LICENSE. The finetuned CLIP-based model is subject to original OPENAI's licensing terms, a MIT-based license: <https://github.com/openai/CLIP/blob/main/LICENSE>.

Acknowledgments: For the Politecnico di Milano's GeoLAB, where most of the present research was physically developed.

Conflicts of Interest: The authors declare no conflicts of interest.

References

1. Hamim, O.F.; Kancharla, S.R.; Ukkusuri, S. Mapping Sidewalks on a Neighborhood Scale from Street View Images. *Environment and Planning B: Urban Analytics and City Science* **2023**, doi:10.1177/23998083231200445.
2. Serna, A.; Marcotegui, B. Urban Accessibility Diagnosis from Mobile Laser Scanning Data. *ISPRS Journal of Photogrammetry and Remote Sensing* **2013**, *84*, 23–32, doi:10.1016/j.isprsjprs.2013.07.001.
3. Vestena, K. de M.; Camboim, S.; Santos, D.R. dos OSM Sidewalkreator: A QGIS Plugin for an Automated Drawing of Sidewalk Networks for OpenStreetMap. *European Journal of Geography* **2023**, doi:10.48088/ejg.k.ves.14.4.066.084.
4. Wood, J. Sidewalk City: Remapping Public Spaces in Ho Chi Minh City. *Geographical Review* **2016**, *108*, 486–488, doi:10.1111/gere.12239.
5. Zhou, Z.; Lin, Y.; Li, Y. Large Language Model Empowered Participatory Urban Planning 2024.
6. Nadkarni, P.M.; Ohno-Machado, L.; Chapman, W.W. Natural Language Processing: An Introduction. *Journal of the American Medical Informatics Association* **2011**, *18*, 544–551, doi:10.1136/amiajnl-2011-000464.
7. Wang, X.; Ji, L.; Yan, K.; Sun, Y.; Song, R. Expanding the Horizons: Exploring Further Steps in Open-Vocabulary Segmentation. In *Pattern recognition and computer vision*; Liu, Q., Wang, H., Ma, Z., Zheng, W., Zha, H., Chen, X., Wang, L., Ji, R., Eds.; Springer Nature Singapore, 2024; pp. 407–419.
8. Eichstaedt, J.C.; Kern, M.L.; Yaden, D.B.; Schwartz, H.A.; Giorgi, S.; Park, G.; Hagan, C.A.; Tobolsky, V.A.; Smith, L.K.; Buffone, A.; et al. Closed- and Open-Vocabulary Approaches to Text Analysis: A Review, Quantitative Comparison, and Recommendations. *Psychological Methods* **2021**, *26*, 398–427, doi:10.1037/met0000349.
9. Zhu, C.; Chen, L. A survey on open-vocabulary detection and segmentation: Past, present, and future 2023.
10. Zareian, A.; Dela Rosa, K.; Hu, D.H.; Chang, S. Open-Vocabulary Object Detection Using Captions. *CoRR* **2020**, doi:https://arxiv.org/abs/2011.10678.
11. Yang, G.; Ye, Z.; Zhang, R.; Huang, K. A Comprehensive Survey of Zero-Shot Image Classification: Methods, Implementation, and Fair Evaluation. *Applied Computing and Intelligence* **2022**, *2*, 1–31.
12. Lampert, C.H.; Nickisch, H.; Harmeling, S. Attribute-Based Classification for Zero-Shot Visual Object Categorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2014**, *36*, 453–465, doi:10.1109/TPAMI.2013.140.
13. Rohrbach, M.; Stark, M.; Schiele, B. Evaluating Knowledge Transfer and Zero-Shot Learning in a Large-Scale Setting. In *CVPR 2011*; IEEE, 2011; pp. 1641–1648.
14. Zhang, D.; Han, J.; Cheng, G.; Yang, M.-H. Weakly Supervised Object Localization and Detection: A Survey. *IEEE Trans. Pattern Anal. Mach. Intell.* **2021**, 1–1, doi:10.1109/TPAMI.2021.3074313.
15. Vo, H.V.; Siméoni, O.; Gidaris, S.; Bursuc, A.; Pérez, P.; Ponce, J. Active Learning Strategies for Weakly-Supervised Object Detection 2022.
16. Blasiis, M.D.; Benedetto, A.; Fiani, M. Mobile Laser Scanning Data for the Evaluation of Pavement Surface Distress. *Remote Sensing* **2020**, *12*, 942, doi:10.3390/rs12060942.

17. Praticò, F.; Vaiana, R. A Study on the Relationship between Mean Texture Depth and Mean Profile Depth of Asphalt Pavements. *Construction and Building Materials* **2015**, *101*, 72–79, doi:10.1016/j.conbuildmat.2015.10.021.
18. Fidalgo, C.D.; Santos, I.M.; Nogueira, C. de A.; Portugal, M.C.S.; Martins, L.M.T. Urban Sidewalks, Dysfunction and Chaos on the Projected Floor. The Search for Accessible Pavements and Sustainable Mobility. In Proceedings of the Proceedings of the 7th International Congress on Scientific Knowledge; 2021.
19. Vaitkus, A.; Andriejauskas, T.; Šernas, O.; Čygas, D.; Laurinavičius, A. Definition of concrete and composite precast concrete pavements texture. *Transport* **2019**, doi:10.3846/transport.2019.10411.
20. Zeng, Z.; Boehm, J. Exploration of an Open Vocabulary Model on Semantic Segmentation for Street Scene Imagery. *ISPRS International Journal of Geo-Information* **2024**, *13*, 153, doi:10.3390/ijgi13050153.
21. Cordts, M.; Omran, M.; Ramos, S.; Rehfeld, T.; Enzweiler, M.; Benenson, R.; Franke, U.; Roth, S.; Schiele, B. The Cityscapes Dataset for Semantic Urban Scene Understanding. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition; 2016; pp. 3213–3223.
22. Geiger, A.; Lenz, P.; Stiller, C.; Urtasun, R. Vision Meets Robotics: The KITTI Dataset. *The International Journal of Robotics Research* **2013**, *32*, 1231–1237.
23. Ros, G.; Sellart, L.; Materzynska, J.; Vazquez, D.; Lopez, A.M. The SYNTHIA Dataset: A Large Collection of Synthetic Images for Semantic Segmentation of Urban Scenes. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition; 2016; pp. 3234–3243.
24. Yu, H.; Yang, Z.; Tan, L.; Wang, Y.; Sun, W.; Sun, M.; Tang, Y. Methods and Datasets on Semantic Segmentation: A Review. *Neurocomputing* **2018**, *304*, 82–103, doi:10.1016/j.neucom.2018.03.037.
25. Hao, S.; Zhou, Y.; Guo, Y. A brief survey on semantic segmentation with deep learning. *Neurocomputing* **2020**, *406*, 302–321, doi:10.1016/j.neucom.2019.11.118.
26. Mo, Y.; Wu, Y.; Yang, X.; Liu, F.; Liao, Y. Review the State-of-the-Art Technologies of Semantic Segmentation Based on Deep Learning. *Neurocomputing* **2022**, *493*, 626–646, doi:10.1016/j.neucom.2022.01.005.
27. Zou, J.; Guo, W.; Wang, F. A Study on Pavement Classification and Recognition Based on VGGNet-16 Transfer Learning. *Electronics* **2023**, *12*, 3370, doi:10.3390/electronics12153370.
28. Zhang, C.; Nateghinia, E.; Miranda-Moreno, L.F.; Sun, L. Pavement Distress Detection Using Convolutional Neural Network (CNN): A Case Study in Montreal, Canada. *International Journal of Transportation Science and Technology* **2022**, *11*, 298–309, doi:10.1016/j.ijtst.2021.04.008.
29. Riid, A.; Lõuk, R.; Pihlak, R.; Tepljakov, A.; Vassiljeva, K. Pavement Distress Detection with Deep Learning Using the Orthoframes Acquired by a Mobile Mapping System. *Applied Sciences* **2019**, *9*, 4829, doi:10.3390/app9224829.
30. Mesquita, R.; Ren, T.I.; Mello, C.; Silva, M. Street Pavement Classification Based on Navigation through Street View Imagery. *AI & SOCIETY* **2022**, doi:10.1007/s00146-022-01520-0.
31. Hosseini, M.; Miranda, F.; Lin, J.; Silva, C.T. CitySurfaces: City-scale semantic segmentation of sidewalk materials. *CoRR* **2022**, 2201.

32. Radford, A.; Kim, J.W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. Learning Transferable Visual Models from Natural Language Supervision 2021.
33. Liu, S.; Zeng, Z.; Ren, T.; Li, F.; Zhang, H.; Yang, J.; Li, C.; Yang, J.; Su, H.; Zhu, J.; et al. Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection 2023.
34. Grinberger, A.Y.; Minghini, M.; Juhász, L.; Yeboah, G.; Mooney, P. OSM Science—The Academic Study of the OpenStreetMap Project, Data, Contributors, Community, and Applications. *IJGI* **2022**, *11*, 230, doi:10.3390/ijgi11040230.
35. Zeng, Y.; Huang, Y.; Zhang, J.; Jie, Z.; Chai, Z.; Wang, L. Investigating Compositional Challenges in Vision-Language Models for Visual Grounding. In Proceedings of the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR); June 2024; pp. 14141–14151.
36. Rajabi, N.; Kosecka, J. Q-GroundCAM: Quantifying Grounding in Vision Language Models via GradCAM 2024.
37. Wang, S.; Kim, D.; Taalimi, A.; Sun, C.; Kuo, W. Learning Visual Grounding from Generative Vision and Language Model 2024.
38. Quarteroni, S.; Dinarelli, M.; Riccardi, G. Ontology-Based Grounding of Spoken Language Understanding. In Proceedings of the 2009 IEEE Workshop on Automatic Speech Recognition & Understanding; IEEE: Moreno, Italy, December 2009; pp. 438–443.
39. Baldazzi, T.; Bellomarini, L.; Ceri, S.; Colombo, A.; Gentili, A.; Sallinger, E. Fine-Tuning Large Enterprise Language Models via Ontological Reasoning 2023.
40. Jullien, M.; Valentino, M.; Freitas, A. Do Transformers Encode a Foundational Ontology? Probing Abstract Classes in Natural Language 2022.
41. FRC CSC RAS / Moscow, Russia; Larionov, D.; RUDN University / Moscow, Russia; Shelmanov, A.; Skoltech / Moscow, Russia; FRC CSC RAS / Moscow, Russia; Chistova, E.; FRC CSC RAS / Moscow, Russia; RUDN University / Moscow, Russia; Smirnov, I.; et al. Semantic Role Labeling with Pretrained Language Models for Known and Unknown Predicates. In Proceedings of the Proceedings - Natural Language Processing in a Deep Learning World; Incoma Ltd., Shoumen, Bulgaria, October 22 2019; pp. 619–628.
42. Smith, M.K.; Welty, C.; McGuinness, D.L. OWL Web Ontology Language Guide 2004.
43. Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A.C.; Lo, W.Y.; et al. Segment anything 2023.
44. Meta Mapillary. *GitHub* **XXXX**.
45. Vestena, K. *GitHub* - kauevestena/deep_pavements_dataset. *GitHub* **XXXX**.
46. Fan, Q.; Tao, X.; Ke, L.; Ye, M.; Zhang, Y.; Wan, P.; Wang, Z.; Tai, Y.-W.; Tang, C.-K. Stable Segment Anything Model 2023.
47. Hetang, C.; Xue, H.; Le, C.; Yue, T.; Wang, W.; He, Y. Segment Anything Model for Road Network Graph Extraction 2024.
48. Son, J.; Jung, H. Teacher–Student Model Using Grounding DINO and You Only Look Once for Multi-Sensor-Based Object Detection. *Applied Sciences* **2024**, *14*, 2232, doi:10.3390/app14062232.

49. Dong, X.; Bao, J.; Zhang, T.; Chen, D.; Gu, S.; Zhang, W.; Yuan, L.; Chen, D.; Wen, F.; Yu, N. CLIP Itself Is a Strong Fine-Tuner: Achieving 85.7% and 88.0% Top-1 Accuracy with ViT-B and ViT-L on ImageNet 2022.
50. Nguyen, T.; Ilharco, G.; Wortsman, M.; Oh, S.; Schmidt, L. Quality Not Quantity: On the Interaction between Dataset Design and Robustness of Clip. *Advances in Neural Information Processing Systems* **2022**, *35*, 21455–21469.
51. Fang, A.; Ilharco, G.; Wortsman, M.; Wan, Y.; Shankar, V.; Dave, A.; Schmidt, L. Data Determines Distributional Robustness in Contrastive Language Image Pre-Training (Clip). In *Proceedings of the International Conference on Machine Learning*; PMLR, 2022; pp. 6216–6234.
52. Tu, W.; Deng, W.; Gedeon, T. A Closer Look at the Robustness of Contrastive Language-Image Pre-Training (Clip). *Advances in Neural Information Processing Systems* **2024**, *36*.
53. Mumuni, F.; Mumuni, A. Segment Anything Model for Automated Image Data Annotation: Empirical Studies Using Text Prompts from Grounding DINO 2024.
54. Kaue-Vestena/Clip-Vit-Base-Patch32-Finetuned-Surface-Materials. Hugging Face. Title of Thesis. Level of Thesis, Degree-Granting University, Location of University, Date of Completion, 2024.
55. Eimer, T.; Lindauer, M.; Raileanu, R. Hyperparameters in Reinforcement Learning and How To Tune Them 2023.
56. Tong, Q.; Liang, G.; Bi, J. Calibrating the Adaptive Learning Rate to Improve Convergence of ADAM. *Neurocomputing* **2022**, *481*, 333–356, doi:10.1016/j.neucom.2022.01.014.
57. Reddi, S.J.; Kale, S.; Kumar, S. On the Convergence of Adam and Beyond 2019.
58. Yu, J.; Wang, Z.; Vasudevan, V.; Yeung, L.; Seyedhosseini, M.; Wu, Y. CoCa: Contrastive Captioners Are Image-Text Foundation Models 2022.
59. Code, P.W. Papers with Code - ImageNet Benchmark (Image Classification. *GitHub* **2024**.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.