

Robust Burned Area Delineation through Multitask Learning ^{*}

Edoardo Arnaudo ^{*1,2[0000-0001-9972-599X]}, Luca Barco ^{*2[0000-0002-9089-9616]},
Matteo Merlo ^{†1[0009-0002-8008-5093]}, and Claudio Rossi ^{2[0000-0001-5038-3597]}

¹ Politecnico di Torino, Corso Duca degli Abruzzi 24, 10129 Torino, Italy
`name.surname@polito.it`, `†s287576@studenti.polito.it`

² Fondazione LINKS, Via Carlo Boggio 61, 10138 Torino, Italy
`name.surname@linksfoundation.com`

Abstract. In recent years, wildfires have posed a significant challenge due to their increasing frequency and severity. For this reason, accurate delineation of burned areas is crucial for environmental monitoring and post-fire assessment. However, traditional approaches relying on binary segmentation models often struggle to achieve robust and accurate results, especially when trained from scratch, due to limited resources and the inherent imbalance of this segmentation task. We propose to address these limitations in two ways: first, we construct an ad-hoc dataset to cope with the limited resources, combining information from Sentinel-2 feeds with Copernicus activations and other data sources. In this dataset, we provide annotations for multiple tasks, including burned area delineation and land cover segmentation. Second, we propose a multitask learning framework that incorporates land cover classification as an auxiliary task to enhance the robustness and performance of the burned area segmentation models. We compare the performance of different models, including UPerNet and SegFormer, demonstrating the effectiveness of our approach in comparison to standard binary segmentation.

Keywords: Remote Sensing · Computer Vision · Semantic Segmentation.

1 Introduction

In recent years, wildfire events have become a recurring major problem, due to their increasing frequency and severity. These events have serious environmental and socio-economic impacts: with potential to cause extensive damage to forests, wildlife habitats, and even human lives. For this reason, understanding and effectively managing wildfires represents a crucial task for first responders and decision makers.

^{*} This work was carried out in the context of the projects SAFERS (H2020, Grant ID. 869353) and OVERWATCH (HEU, Grant ID. 101082320).

^{*} Equal contribution.

Accurate and reliable delineation of burned areas is therefore essential for various applications, including environmental monitoring and post-fire assessment. Traditional approaches for burned area delineation often rely on binary segmentation models trained from scratch. However, these models may struggle to achieve accurate and robust results, due to the limitations of the underlying data. First, resources specifically tailored for this task remain particularly scarce, often lacking large and diverse datasets for an effective training. Second, burned area segmentation is an inherently unbalanced problem, as the extent of burned areas is often significantly smaller compared to non-burned regions in the input imagery. This imbalance usually hinders the generalization abilities of the models, when applied in different scenarios.

Furthermore, existing datasets used for burned area delineation are often lacking in terms of surface covered [8] or diversity [17]. These shortcomings may hinder the ability of the models to generalize effectively, underlining the need for more comprehensive and varied data sources to enhance model performance and real-world applicability.

To address these limitations, we first construct an ad-hoc dataset, specifically tailored for the task of burned area segmentation, cross-referencing information from the Copernicus European Monitoring System (EMS) with Sentinel-2 feeds and other relevant sources. This dataset provides a comprehensive set of samples, with a focus on the European soil, including annotations for two different tasks: burned area delineation, and land cover segmentation. Second, we propose a multitask learning framework that leverages land cover classification as an auxiliary task. By incorporating this information into the learning process, we aim to improve the robustness and performance of the model on the burned area segmentation task. We compare the performance of different models, including UPerNet [23] and SegFormer [24], demonstrating the effectiveness of our approach against a classic binary segmentation training in several configurations, including with and without bootstrap from pretrained weights. Dataset and code related to this work are available at github.com/links-ads/burned-area-seg.

The remainder of this paper is structured as follows. Section 2 reviews related works, Section 3 describes the construction of the multitask dataset for burned area delineation, while Section 4 presents the proposed multitask learning framework and its components. Section 5 details the experimental setup and discusses the obtained results. Lastly, Section 6 concludes the manuscript, suggesting potential future directions.

2 Related Works

2.1 Aerial Semantic Segmentation

Considering remote sensing and aerial images, semantic segmentation plays a crucial role in various applications, including urban planning [13], land cover monitoring [2], and crisis management [5]. Existing semantic segmentation methods typically rely on convolutional encoder-decoder architectures (CNNs), with

different variants to capture both the global context and the finer details of the scene. Approaches such as Fully Convolutional Networks (FCN) [12] and U-Nets [19], make use of bottleneck components to encode pixel information into semantically meaningful vectors, coupled with skip connections to integrate lower-level features. Other solutions involve multiscale feature extraction and fusion, such as DeepLab [6] and PSPNet [25], where inputs are processed with varying kernel sizes and dilations to capture local and global context at once. Subsequent variants often combine these concepts to provide more robust features [6,23]. Segmenting aerial images introduces several specific challenges that often require tailored solutions. Unlike other settings, satellite data often provides multiple spectra beyond the visible bands. These can be integrated in multiple ways, such as additional input channels [20] or using ad-hoc encoders for feature fusion [21]. Moreover, aerial images they are often denser, containing several entities against complex backgrounds, with wider spatial relationships. To address this, attention components are commonly employed to better model long-distance similarities among pixels [16]. Transformer architectures and their segmentation variants [24] thus become a natural choice in this case, given its inherent ability at extracting long-range relations.

2.2 Burned Area Delineation

Over the years, numerous techniques have been proposed to delineate burned areas from remote sensing data. Standard approaches make use of several spectral indices, to discern burned soil from unburned areas using a combination of multiple bands. The *de facto* standard is represented by the Normalized Burn Ratio (NBR) [7] and the difference NBR (dNBR) [14], which are often used in combination with other indices [9]. Other variants have been developed to better adapt to specific satellite feeds, such as the Burned Area Index for Sentinel-2 (BAIS) [10]. However, these approaches are usually noisy and require further manual processing to produce clean results. In some cases, such as the dNBR, a pre-wildfire image is also needed to compare the same regions before and after the event. In the last decades, machine and deep learning techniques obtained promising results, reducing the manual effort while obtaining more robust tools. Supervised classification algorithms, such as Support Vector Machines (SVM) and Random Forests (RF), have been widely employed for burned area mapping [18,11]. Acting on a per-pixel basis, these approaches remain effective on lower resolution feeds such as MODIS, however their lack of contextual information may result in suboptimal results on higher resolution feeds such as Sentinel-2 [11]. Recently, convolutional networks have been successfully employed to produce robust results on this task, especially considering post-wildfire images only. U-net segmentation architectures [15,11] represent the standard approach, however Transformer-based architectures have demonstrated their effectiveness in several remote sensing scenarios [20], including burned area segmentation [5].

3 Dataset

To carry out this work, a crucial initial step involved the construction of a custom dataset specifically tailored for multitask learning, focusing on wildfire events. Expanding on similar works in this field [8], our dataset contains 171 fire events derived from Copernicus EMS³. For each Area Of Interest of the event, we provide (i) the Sentinel-2 satellite imagery, (ii) a burned area annotation, derived from EMS (iii) a land cover map, derived from ESA WorldCover, and (iv) a cloud mask computed on the remote sensing input.

3.1 Data sources

We gather all the available large wildfire events in recent years from the catalog of the Copernicus EMS, an integral part of the Copernicus program launched by the European Union. Within this open service, the Rapid Mapping module plays a crucial role by providing a curated set of Areas of Interest (AoI) associated with each event, where each crisis event has been carefully analyzed, and its delineation has been manually generated by a team of experts. Every AoI may provide three distinct manual annotations, named products: the First Estimate Product (FEP), the Delineation Product, and the Grading Product. The FEP consists of preliminary information about the affected territory and event, facilitating initial emergency response efforts. On the other hand, the delineation and grading products offer a more accurate and comprehensive label regarding the extent of the event and the assessment of damages. Following the geographical coordinates and the time of the event associated with each AoI, we download the corresponding satellite images from the Sentinel-2 mission, that serve as input to the deep learning algorithms. Sentinel-2 captures data across 12 spectral bands with varying resolutions, ranging from 10 to 60 meters. In this study, we focus on the L2A product, which transforms the reflectance into Bottom-of-Atmosphere (BoA) values through atmospheric correction.

In addition to the satellite imagery and the burned area delineation maps, we also incorporate land cover data on the same area as an auxiliary target, exploiting the ESA World Cover dataset [1]. This resource offers annual maps for the years 2020 and 2021, featuring 11 distinct classes, including trees, shrubland, grassland, built-up areas, areas with sparse vegetation, water bodies and other surfaces. By integrating a more generic land cover information, we aim to enrich the semantic segmentation process with a broader understanding of the landscape dynamics, enhancing the robustness and contextual accuracy of our model.

3.2 Data preparation

We download each EMS activation available, with the aim of maximizing the amount of samples with valid input image and corresponding ground truth labels. For each fire event we gather several details, including the event date,

³ <https://emergency.copernicus.eu/>

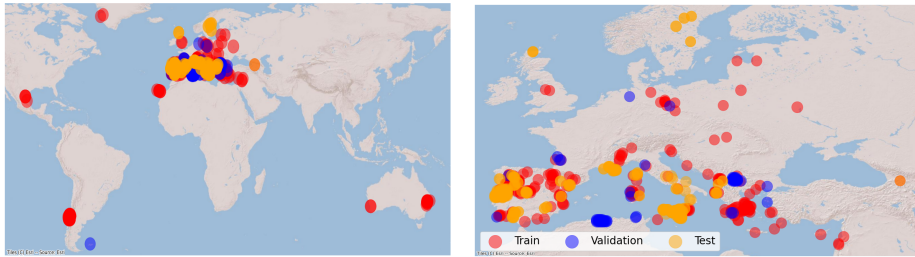


Fig. 1. Distribution of fire events contained in our dataset divided into train (red), validation (blue), and test (orange), on a worldwide scale (left) and at European level (right).

geographical coordinates defining the bounding box of the affected area, and corresponding delineation and grading maps. However, we note that there may be cases where the delineation, the grading, or both maps are unavailable. In such instances, given its higher quality, we maintain only those areas with a valid grading, and we generate the corresponding delineation map performing a standard binarization over the burn severity values. Starting from the remaining processed activations, we retrieve the corresponding post-fire Sentinel-2 images, exploiting the SentinelHub services⁴. Given the input requirements of the models, we force each image to have a minimum dimension of 512 pixels on each side, expanding the smaller regions until this requirement is satisfied for every AoI. At the same time, we split areas larger than $2,500 \times 2,500$ pixels in multiple subsections for practical use. We sample and rasterize each image with a resolution of 10m per pixel, the maximum provided by Sentinel-2, upscaling the lower resolution bands with nearest neighbor interpolation. To maximize the number of clear images, without smoke or large clouds, we consider a time frame of up to 30 days following the reported event date, selecting the satellite acquisition with the least cloud coverage. Despite these precautions, it is not uncommon to observe clouds in the final image samples: for this reason, we further process the images using a cloud segmentation model, derived from CloudSen12 [3], generating a validity map. This additional mask is then applied during training, excluding every pixel covered by clouds from the loss computation. For the corresponding land cover maps, we retrieve the required raster layers from the ESA World Cover database, available via Microsoft Planetary Computer⁵. No further processing is applied to the labels, except for a direct remapping from the original ESA taxonomy to a contiguous list of categories indexed from 0. A value of 255 is further assigned to pixels missing their specific category.

The final dataset comprises a collection of 433 samples, spanning from 2017 to the first months of 2023. The events are predominantly concentrated in Europe, with select events occurring in Australia and on the American continent. Given

⁴ <https://www.sentinel-hub.com/>

⁵ <https://planetarycomputer.microsoft.com/>

the same source, our collection effectively represents an extension of previous datasets [8]. For this reason, we dedicate every activation already present in previous works to testing purposes, training on the remaining events. This allows for easier comparisons with prior results, while also serving as a benchmark for assessing the generalizability and performance of our proposed approach.

4 Methodology

4.1 Problem statement

The problem at hand involves developing a multitask learning framework for burned area delineation, exploiting land cover classification as an auxiliary target to guide the training. We have access to a delineation map (y_D) and a land cover map (y_{LC}) as ground truth labels. We employ models composed of a single encoder and a single decoder with two classification heads, namely h_D and h_{LC} . The objective is to simultaneously train the model f_θ , with parameters θ , to predict accurate burned area delineations (\hat{y}_D), while also training on land cover classification (\hat{y}_{LC}) using the shared representations ϕ_θ . The shared architecture with two standard classification heads enables the model to learn from both tasks jointly.

4.2 Framework and models

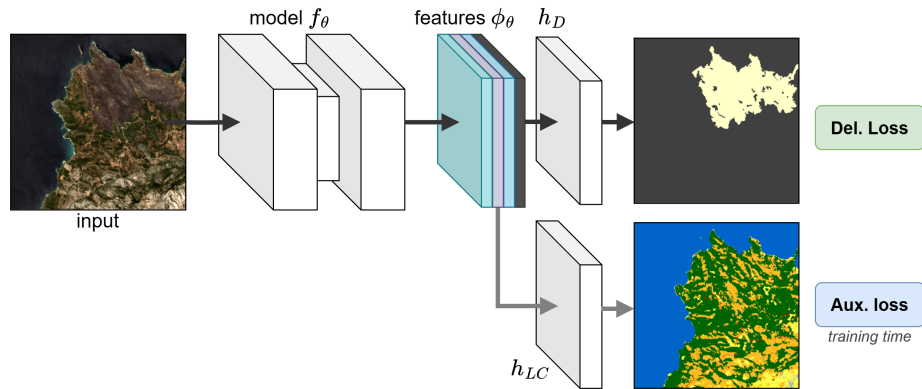


Fig. 2. Multitask learning framework: the decoder features are shared with the auxiliary head h_{LC} for joint training. The auxiliary head is dropped at test time.

Our approach is shown in Fig. 2. To train the full model f_θ we simultaneously predict burned area delineation and land cover segmentation using the shared representations from the decoder stage ϕ_θ . These enable the model to capture and leverage common patterns and features between the two tasks, which may

help in improving the segmentation outcome. Throughout the training process, we employ a standard Cross Entropy loss, in its binary and multi-class variants respectively. The gradients derived from both tasks are jointly propagated back to update the model’s parameters. At test time, we drop the auxiliary head, focusing only on the burned area delineation performance, through standard binary segmentation.

To provide a comprehensive overview and compare standard convolutional networks with vision transformers, we explore three different architectures: two UPerNet [23] variants, using a Residual Network (ResNet) and a Vision Transformer (ViT) as encoders respectively, and SegFormer [24]. Thanks to its unified perceptual parsing structure, UPerNet provides the flexibility to use both standard CNNs, and recent transformer-based solutions. This allows for a better comparison between the two architectures. On the other hand, SegFormer represents an alternative end-to-end solution which demonstrated its effectiveness on aerial tasks, including burned area delineation [5,20].

5 Experiments

5.1 Implementation details

As mentioned in Section 3, we train our models on the subset of activations that are not present in previous datasets [8], considering the remaining ones as our test set. We further extract a 10% of activations from our training for validation purposes, obtaining a total of 129 wildfire events in training, 15 in validation, and 27 for testing purposes. To cope with the varying image dimensions, we implement a random sampling strategy that extracts square crops of 512×512 pixels from random image sections at runtime during training. For validation and testing, we adopt instead a sequential sampling strategy with overlapping tiles, reconstructing the original inputs by means of a smooth blending using splines. We consider two groups of experiments: first, we only focus on burned area delineation, as a single training task. Second, we conduct a multitask training, using both delineation and land cover maps. In the latter case, we further mask out the burned pixels from the annotation, to avoid inconsistent labels. For both scenarios, we train three architectures: UperNet with two different encoders (i.e., ResNet-50 and ViT-S) and SegFormer with MiT-B3 as encoder. Moreover, we investigate the impact of using pretrained weights on the backbones in both configurations. We exploit pretrained weights derived from large-scale pretraining on SSL4EO-S12 [22] for ResNet and ViT, in the RN50 and ViT-S variants, while we adopt weights pretrained on ImageNet [24] for SegFormer for lack of better options. We base our code on the *mmsegmentation*⁶ library, adapting the model inputs to adjust for the additional channels. In every experiment, we train on a single NVIDIA A100 GPU for 30 epochs, using a batch size of 32 tiles, AdamW as optimizer with a learning rate of $1e-4$, and a Cross Entropy loss for both tasks, in binary and multi-class versions respectively. Following similar works [5,20], we

⁶ <https://github.com/open-mmlab/msegmentation>

adopt macro-averaged F1 score and Intersection over Union (IoU) as evaluation metrics in every configuration.

5.2 Results

Setting	Model	From scratch		Pretrained	
		F1	IoU	F1	IoU
STL	SegFormer (MiT-B3)	89.01± 1.39	80.22± 2.25	90.79± 0.46	83.13± 0.78
	UPerNet (RN50)	82.33± 9.17	70.94± 12.63	91.27± 0.08	83.95± 0.13
	UPerNet (ViT-S)	87.65± 2.01	78.08± 3.17	89.20± 1.29	80.53± 2.09
MTL	SegFormer (MiT-B3)	90.94± 0.17	83.38± 0.29	90.91± 0.28	83.34± 0.47
	UPerNet (RN50)	89.82± 1.76	81.57± 2.87	91.86± 0.30	84.94± 0.51
	UPerNet (ViT-S)	89.76± 0.15	81.43± 0.25	90.69± 0.58	82.98± 0.97

Table 1. Experimental results in single (STL) and multitask (MTL) training, comparing models trained from scratch, or using pretrained encoders.

Setting	Model	Training time (1 Ep.)	Param. (M)
STL	SegFormer (MiT-B3)	3h28m	44,6
	UPerNet (RN50)	3h20m	64,1
	UPerNet (ViT-S)	3h20m	57,9
MTL	SegFormer (MiT-B3)	3h40m (+12m)	44,6
	UPerNet (RN50)	3h50m (+30m)	64,1
	UPerNet (ViT-S)	3h30m (+10m)	57,9

Table 2. Analysis of the computational costs in terms of training time over one epoch, as average of three epochs, and total network parameters. While the training time increases by a small margin, the parameter increase is effectively negligible given the shared encoder-decoder structure.

We conduct single-task (STL) and multitask (MTL) training experiments using both pretrained and non-pretrained weights. For each configuration, we perform three separate runs with different seeds, reporting the results in Table 5.2 as average scores with their corresponding standard deviation. Focusing on the experiments conducted without pretrained weights, the multitask setting consistently achieves superior performance and lower standard deviation compared to the single task setting. Except for the SegFormer, that reports the highest scores in both variants, the multitask approach exhibits a noticeable improvement of +3.85 in terms of F1 score, or +5.71 in terms IoU, averaged across every model. Furthermore, we note that the results are way more stable in the multitask configuration, where the standard deviation decreases by -3.51 (F1) and -4.88 (IoU). This is also shown in Figure 3, where the latter produce more reliable segmentation maps. Considering the experiments with pretrained

weights, the disparity between single and multitask performance is no longer apparent, with higher and more stable scores even in the single task setup. This is expected, as large-scale pretraining has been proven to be effective in several contexts [22]. Nevertheless, multitask training still yields an average overall improvement of +0.73 in F1 score and +1.21 in IoU, regardless of the underlying architecture. Lastly, comparing training from scratch and using pretrained weights, the latter always exhibit higher performances, even more so in multitask configuration. Specifically, when comparing the top performing models from both tables (i.e., Segformer in the first case, UPerNet-RN50 in the second case), the single task setting achieves +2.24 in F1 score and +3.72 in IoU, whereas in multitask achieves +0.92 in F1 score and +1.56 in IoU. Overall, the results demonstrate the validity of the multitask strategy, exhibiting increased performance robustness, comparable to or even surpassing pretrained solutions in certain instances.

In Table 5.2 we also compare the computational costs of the STL approach compared to the MTL solution. We observe that the MTL versions, despite including an additional segmentation head, exhibit only a modest increase in training speed compared to their single task learning (STL) counterparts, with a marginal difference of 20 seconds in training speed. Moreover, while MTL models do incur a slight increase in memory usage, this increment remains negligible and does not substantially impact the feasibility of implementation. This is expected, since the MTL setting only effectively adds the parameters of a single pixel classification head, which boils down to a 1×1 convolutional layer with $|\phi_\theta|$ feature channels as input, and 11 categories as output. Moreover, we note that during inference the auxiliary head is omitted, effectively eliminating any computational overhead associated with the auxiliary task.

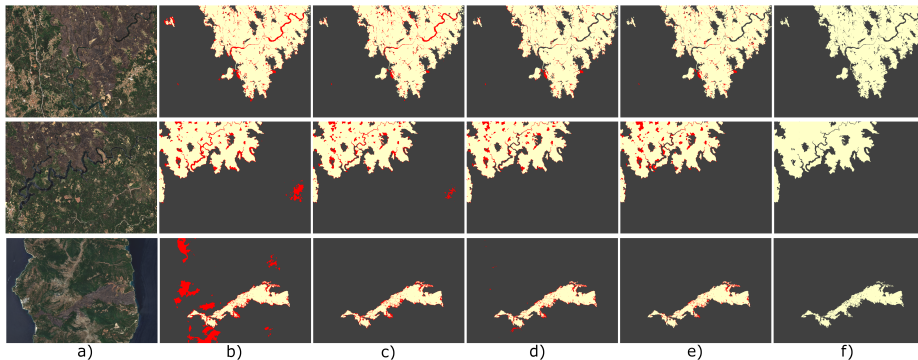


Fig. 3. Qualitative examples derived from UPerNet-RN50: a) Sentinel-2 input; b) Single task and c) Multitask from scratch; d) Single Task and e) Multitask with pretrained weights; f) ground truth. Red pixels represent prediction errors.

6 Conclusion

In this work, we propose a multitask approach for burned area delineation, exploiting land cover classification as an auxiliary target. Results show that the devised solutions yield more stable and robust performances, comparable to or even superior to pretrained solutions.

Multitask learning offers promising results, especially in absence of pretrained solutions. Despite the robust performance, the current multitask approach presents some limitations: first, the performance of the models heavily rely on the quality of the annotations of both tasks. Second, the improved robustness comes at the cost of additional computational complexity, which may limit the scalability. Moreover, we recognize the need to delve deeper into the impact of task characteristics and explore a wider array of auxiliary tasks for a more comprehensive multi-task learning approach, including for instance multiple training objectives to further enhance scalability and generalization capabilities of the models.

Future studies may therefore focus on improving the multitask capabilities by integrating multiple heterogeneous tasks at the same time [4], or may consider more computationally demanding approaches such as large-scale self-supervised learning, to generate better pretrained solutions and thus translate these downstream tasks in simpler and faster fine-tuning objectives.

References

1. Agency, E.S.: Esa world cover 2020 (2020), https://worldcover2021.esa.int/data/docs/WorldCover_PUM_V2.0.pdf
2. Arnaudo, E., Tavera, A., Masone, C., Dominici, F., Caputo, B.: Hierarchical instance mixing across domains in aerial segmentation. *IEEE Access* **11**, 13324–13333 (2023)
3. Aybar, C., Ysuhuaylas, L., Loja, J., Gonzales, K., Herrera, F., Bautista, L., Yali, R., Flores, A., Diaz, L., Cuenca, N., et al.: Cloudsen12, a global dataset for semantic understanding of cloud and cloud shadow in sentinel-2. *Scientific data* **9**(1), 782 (2022)
4. Bastani, F., Wolters, P., Gupta, R., Ferdinando, J., Kembhavi, A.: Satlas: A large-scale, multi-task dataset for remote sensing image understanding. *arXiv preprint arXiv:2211.15660* (2022)
5. Cambrin, D.R., Colomba, L., Garza, P.: Vision transformers for burned area delineation. In: *Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases* (2022)
6. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587* (2017)
7. Cocke, A.E., Fulé, P.Z., Crouse, J.E.: Comparison of burn severity assessments using differenced normalized burn ratio and ground data. *International Journal of Wildland Fire* **14**(2), 189–198 (2005)
8. Colomba, L., Farasin, A., Monaco, S., Greco, S., Garza, P., Apiletti, D., Baralis, E., Cerquitelli, T.: A dataset for burned area delineation and severity estimation from satellite imagery. In: *Proceedings of the 31st*

- ACM International Conference on Information & Knowledge Management. p. 3893–3897. CIKM '22, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3511808.3557528>, <https://doi.org/10.1145/3511808.3557528>
9. Escuin, S., Navarro, R., Fernández, P.: Fire severity assessment by using nbr (normalized burn ratio) and ndvi (normalized difference vegetation index) derived from landsat tm/etm images. *International Journal of Remote Sensing* **29**(4), 1053–1073 (2008)
 10. Filippini, F.: Bais2: Burned area index for sentinel-2. In: *Proceedings*. vol. 2, p. 364. MDPI (2018)
 11. Knopp, L., Wieland, M., Rättich, M., Martinis, S.: A deep learning approach for burned area segmentation with sentinel-2 data. *Remote Sensing* **12**(15), 2422 (2020)
 12. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 3431–3440 (2015)
 13. Mahmud, M.N., Osman, M.K., Ismail, A.P., Ahmad, F., Ahmad, K.A., Ibrahim, A.: Road image segmentation using unmanned aerial vehicle images and deeplab v3+ semantic segmentation model. In: *2021 11th IEEE International Conference on Control System, Computing and Engineering (ICCSCE)*. pp. 176–181. IEEE (2021)
 14. Miller, J.D., Thode, A.E.: Quantifying burn severity in a heterogeneous landscape with a relative version of the delta normalized burn ratio (dnbr). *Remote Sensing of Environment* **109**(1), 66–80 (2007)
 15. Monaco, S., Greco, S., Farasin, A., Colomba, L., Apiletti, D., Garza, P., Cerquitelli, T., Baralis, E.: Attention to fires: Multi-channel deep learning models for wildfire severity prediction. *Applied Sciences* **11**(22), 11060 (2021)
 16. Niu, R., Sun, X., Tian, Y., Diao, W., Chen, K., Fu, K.: Hybrid multiple attention network for semantic segmentation in aerial images. *IEEE Transactions on Geoscience and Remote Sensing* **60**, 1–18 (2021)
 17. Prabowo, Y., Sakti, A.D., Pradono, K.A., Amriyah, Q., Rasyidy, F.H., Bengkulah, I., Ulfa, K., Candra, D.S., Imdad, M.T., Ali, S.: Deep learning dataset for estimating burned areas: case study, indonesia. *Data* **7**(6), 78 (2022)
 18. Ramo, R., Chuvieco, E.: Developing a random forest algorithm for modis global burned area classification. *Remote Sensing* **9**(11), 1193 (2017)
 19. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*. pp. 234–241. Springer (2015)
 20. Tavera, A., Arnaudo, E., Masone, C., Caputo, B.: Augmentation invariance and adaptive sampling in semantic segmentation of agricultural aerial images. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 1656–1665 (2022)
 21. Valada, A., Vertens, J., Dhall, A., Burgard, W.: Adapnet: Adaptive semantic segmentation in adverse environmental conditions. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*. pp. 4644–4651. IEEE (2017)
 22. Wang, Y., Braham, N.A.A., Xiong, Z., Liu, C., Albrecht, C.M., Zhu, X.X.: Ssl4eos12: A large-scale multi-modal, multi-temporal dataset for self-supervised learning in earth observation (2023)

23. Xiao, T., Liu, Y., Zhou, B., Jiang, Y., Sun, J.: Unified perceptual parsing for scene understanding. In: Proceedings of the European conference on computer vision (ECCV). pp. 418–434 (2018)
24. Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., Luo, P.: Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems* **34**, 12077–12090 (2021)
25. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2881–2890 (2017)