

---

# COMPARING A COMPOSITE MODEL VERSUS CHAINED MODELS TO LOCATE A NEAREST VISUAL OBJECT

---

**Antoine Le Borgne**  
Orange Innovation  
Lannion  
antoine.leborgne@orange.com

**Xavier Marjou**  
Orange Innovation  
Lannion  
xavier.marjou@orange.com

**Fanny Parzysz**  
Orange Innovation  
Lannion  
fanny.parzysz@orange.com

**Tayeb Lemlouma**  
IRISA  
Lannion  
tayeb.lemlouma@irisa.fr

## ABSTRACT

Extracting information from geographic images and text is crucial for autonomous vehicles to determine in advance the best cell stations to connect to along their future path. Multiple artificial neural network models can address this challenge; however, there is no definitive guidance on the selection of an appropriate model for such use cases. Therefore, we experimented two architectures to solve such a task: a first architecture with chained models where each model in the chain addresses a sub-task of the task; and a second architecture with a single model that addresses the whole task. Our results showed that these two architectures achieved the same level performance with a root mean square error (RMSE) of 0.055 and 0.056; The findings further revealed that when the task can be decomposed into sub-tasks, the chain architecture exhibits a twelve-fold increase in training speed compared to the composite model. Nevertheless, the composite model significantly alleviates the burden of data labeling.

## 1 Introduction

Extracting information from geographic images and text is crucial to predict connectivity for autonomous vehicles. Autonomous vehicles require excellent cellular connectivity to interact with other connected vehicle, cars analyze images to locate themselves in their environment or even unmanned aerial vehicles (UAVs) for use-case like surveillance of asset over long geographical distances. Consequently, most operators want to assess the radio reception of their autonomous vehicle along their planned trajectories. Various connectivity-related features can be pertinent in predicting the cellular connectivity of a given area: the number of base stations, their geographic position and azimuth of the antenna, the topology of the area, the characteristics of the area (rural, urban) or the population density. Typically, these inputs come from tabular data and map images, thus requiring multi-modal models.

With recent advances in machine learning to perform multi-modal tasks, recent studies show that in relatively near future, Large Language models should be the most suitable model for multi-modal tasks [1]. Large Language Models (LLMs) are neural networks with billions of parameters that are trained on vast amounts of unlabelled text using self-supervised learning. LLMs have emerged as a turning point in natural language processing research, shifting the focus away from specialized supervised models for specific tasks. Recently, GPT4.0, a generative pre-trained Transformer (GPT) based model [2], was announced to be a "large multi-modal model (accepting image and text inputs, emitting text outputs)"[3]. Although image support is not yet officially open to users, in its current state, GPT4.0 seems capable of roughly recognizing image content, but not capable of specific tasks such as measuring distance, as shown by an example in appendix. To be fair with GPT4.0, additional testing would be needed given that we only experimented it in one-shot or few-shot learning, not a full learning as with our methods.

Another LLM that has reached the state of the art is LLaMA presented by [4]. And in recent work, [1] presents an adapter that can be tuned in a few hours to enable LLaMA to support multi-modal input as images, it may be interesting for the future to look at this work.

However, in our work we focused on extracting information only from images. A first alternative is to divide this task into multiple sub-task and chain a sequence of models, one model dedicated to each sub-task. But for each of the sub-tasks, this requires significant labeling work that we want to avoid. A second alternative is to use a single model that we call composite model, to entirely solve the task. To our knowledge, no study has compared whether it was better to use this composite model or chained models composed of several models specific to each sub-task.

This article focuses on the comparison of two models that aim to predict the  $(x, y)$  coordinates of the base station closest to a drone based on an input image. This prediction task involves two essential sub-tasks: information extraction from the image and subsequent analysis to determine the  $(x, y)$  coordinates. Deep learning models such as Convolutional Neural Networks (CNNs) and Transformers have demonstrated their efficacy in these areas, making them suitable choices for accomplishing these tasks effectively.

Convolutional Neural Networks (CNNs) have become increasingly popular for regression tasks in recent years [5] [6] [7]. Their ability to learn complex feature representations from image data has made them a natural choice for tasks such as image processing. One common approach is to use pre-trained CNNs as feature extractors and feed the extracted features into a linear regression model. This has been successful in applications such as predicting air quality from satellite images [5], estimating forest biomass from UAV-based LiDAR data [6], or detecting COVID-19 from x-ray scans [7]. Another approach is to train CNNs end-to-end for regression tasks, which allows the model to learn both the feature representation and the regression function simultaneously, potentially leading to better performance. For example, [8] proposed an end-to-end CNN model for estimating the height of buildings from aerial images, while [9] used an end-to-end CNN for predicting crop yield from UAV-based multispectral images.

Although CNNs make it possible to carry out advanced classification or regression operations from input images, certain tasks require more complex models, such as Transformer models [10].

Transformer models were first introduced by [10] in the context of machine translation and have been widely used in recent years for various NLP tasks due to their ability to model long-term dependencies and capture contextual relationships. One of the most popular applications of transformers is in language modeling, where the goal is to predict the next word in a sequence given the previous words. It has also been applied to a wide range of other NLP tasks, including text classification, sentiment analysis and question answering. [11] and [12] present transformers models to perform on these tasks.

With the images we seek to analyze, a CNN is sufficient to extract information from the image [13]. However, for our future work we will need models that are flexible and can potentially support multiple inputs. This is why we focused on advanced models capable of analyzing images and incorporating a transformer, a type of machine learning model well known for its flexibility and its ability to take different types of data as input (e.g. images and text). Composite models, also known as hybrid models or end-to-end models are currently used in different fields such as medical field, to analyze medical images or in the field of drones to perform object recognition and classification from UAVs images. Recently google present [14] using multi-task end-to-end transformer model on the task of conversational recommendations, which aim to provide recommendations based on a user's explicit preferences expressed in dialogue.

Several studies have explored the integration of CNNs or transformer models for image classification and object detection tasks. For example, the Vision Transformer model proposed by [15] uses a transformer-based architecture to process image patches and achieve state-of-the-art performance on image classification tasks. Similarly, the Detection Transformer proposed by [16] uses a CNN to extract image features and predict object class in images. There is also the Hybrid CNN-Transformer model proposed by [17] that combines a CNN and a transformer model for object detection tasks and achieves competitive results on benchmark datasets.

The remainder of this article is organized as follows. Section 2 describes the compared models and Section 3 present the experiments. The results and our discussion are in in Section 4. Finally, Section 5 concludes this study.

## 2 Proposed models

This section introduces the models under comparison, namely the chained models, the composite model, and the standalone CNN as a baseline. We outline the specific characteristics of the input data and elaborate on the operational principles of each model.

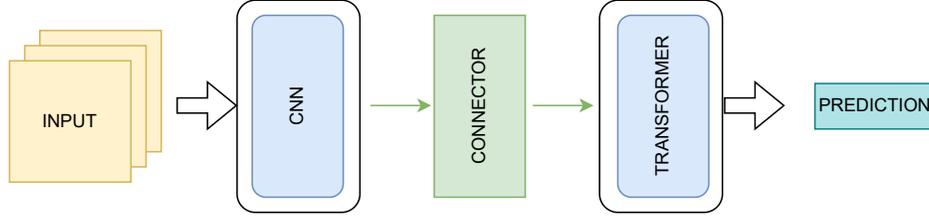


Figure 1: Chained models architecture

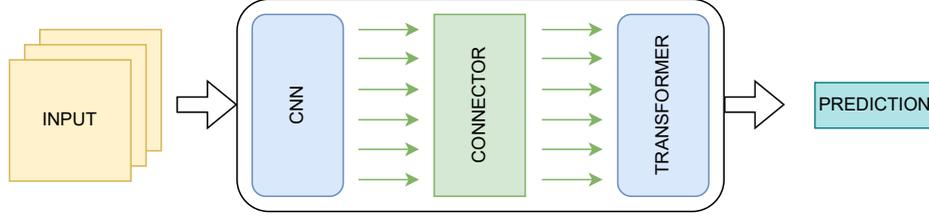


Figure 2: Composite model architecture

## 2.1 Chained models

Chained models, as depicted in Figure 1, are composed of two independently trained models, a CNN and a Transformer. Given that a transformer operates on tokenized inputs, we also add a connector between these models whose role is to sequence and tokenize the CNN outputs to feed into the Transformer. The tokenization is performed character by character with a vocabulary composed of 14 characters: [0,1,2,3,4,5,6,7,8,9, -,.,e,:].

First, each input image of size  $H \times W \times C$  ( $H$ ,  $W$ , and  $C$  respectively stand for height, weight, and channel number) is input into the CNN to produce  $2 \times B + 2 \times U$  output values ( $B$  and  $U$  respectively represent the number of base stations and the number of UAVs in the image). Each of these output values represents an 'x' or 'y' coordinate of the different elements present in each image (the squares and the triangle). They are tokenized to be put as input to the transformer.

The transformer model is trained with the  $(x, y)$  coordinates of the squares and the triangle as inputs and with the  $(x, y)$  coordinates of the square closest to the triangle as the label. The goal is to predict the final  $(x, y)$  value representing the coordinates of the target base station. By training the CNN with the images as input and the  $(x, y)$  coordinates of the squares and the triangle as labels, and subsequently training the transformer with the  $(x, y)$  coordinates of the squares and the triangle as inputs and the position of the closest square as the label, the model learns to generate accurate predictions of the target base station's coordinates.

## 2.2 Composite model

The composite model, as depicted in Figure 2, is quite similar to the chained one but is trained end-to-end.

In this composite model, we use input images of size  $H \times W \times C$ . The CNN and the Transformer are components of the same model and are trained together. Unlike the previous model, the number of CNN outputs can be customized as a hyper-parameter, providing flexibility. This integration of the CNN and Transformer creates a black box-like model, where specific parameters can be modified to enhance performance.

During training, the model is fed with samples comprising an image as the input feature and the  $(x, y)$  coordinates of the square closest to the triangle as the label. The objective is to train the model to accurately predict the coordinates of the target base station. By combining the capabilities of the CNN and the Transformer in a single model and optimizing the parameters, we aim to improve the model's performance.

## 2.3 CNN

We used a simplified version of VGG16 [18] as shown in Figure 3. Our CNN is composed of four convolutional layers with maximum pooling. The input images are of size  $H \times W \times C$ , and the number of CNN outputs  $N$ , is set to 2 to compare our result to the ground truth as a coordinate point.

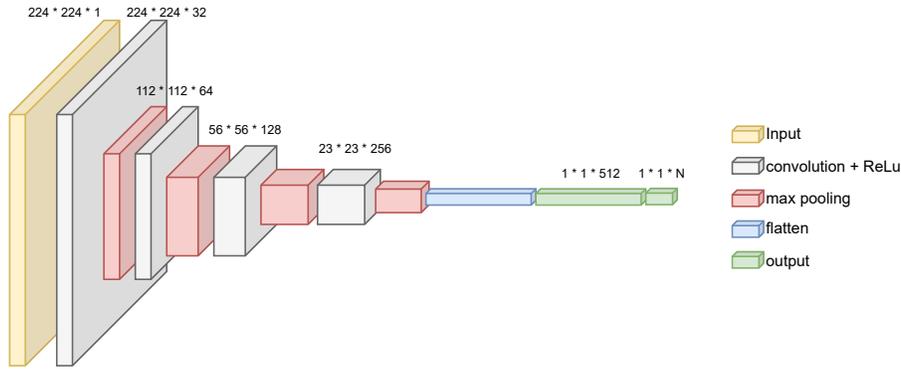


Figure 3: detail of the CNN of our models

### 3 Experiment

This section presents how we optimized our models and how we evaluated them. We specify the hyper-parameters that have been the most influential as well as their optimal values for our use case. We also detail the dataset we use and what composes it.

#### 3.1 Dataset

As we have chosen a simplify task, the dataset is also simplified. The dataset we used is entirely synthetic and specific to this problem. This allowed us to create a large amount of data that can be used to train both a CNN and a transformer. The dataset is composed of 1,000,000 images of size  $224 \times 224$  in gray-scale and the labels associated with it. Some examples can be found in the appendix. The labels associated with each image contain location features that can be used to train our models. Each sample contains a relative path to the relevant image, the  $(x, y)$  coordinates of a first square, the  $(x, y)$  coordinates of a second square, the  $(x, y)$  coordinates of a triangle, and finally, as the ground-truth  $(x, y)$  coordinates of the nearest square to the triangle.

Classically for each training, we separated the dataset into two parts: we kept 1000 images for evaluation and used to the other images to train our models. We also created an additional dataset of 1000 images to test our models.

#### 3.2 Metrics

To evaluate the efficiency of our models, we measured the root mean square error (RMSE) between the point predicted by our model and the ground-truth available in the test dataset. We measured both the average RMSE and the standard deviation to obtain an accuracy metric. In our experiment, the RMSE is between 0 and  $\sqrt{2}$ , a lower average RMSE indicates a better fit between the model and the data.

#### 3.3 Hyper-parameters

To fairly compare the models, we used the same hardware, we used a single Nvidia Geforce RTX 2080 TI. And we chose to tune hyper-parameters that are well known in machine learning to be important, as well as hyper-parameters that have a direct influence on our model.

We modified the learning rate of our model and the size of the dataset to train our models. Two other important hyper-parameters are the number of CNN outputs and the number of decimal places of these output values. They represent the verbosity of the CNN, the amount of information the CNN transfers to the transformer. And the last parameter we changed was the connector function whose role is to organize the CNN information before tokenization, which is essential to make the CNN output understandable by the transformer.

To optimize our models, we employed the Ray Tune library to simultaneously test various parameter combinations. The MLflow library played a crucial role in providing insights into the model's performance, allowing us to assess the impact of different parameter settings and facilitate comparisons between models. This approach allowed us to make informed decisions in selecting the most effective hyper-parameters to maximize the performance of our models.

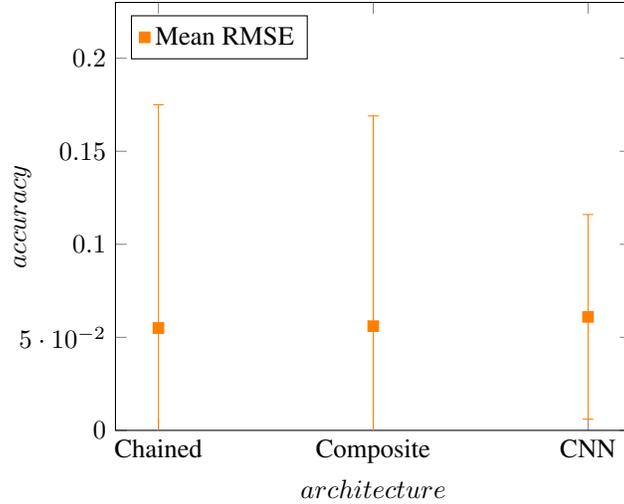


Figure 4: RMSE of the tested models

For the chained models hyper-parameters, the number of values and decimal places in the CNN output were determined by the dataset used to construct the CNN. To compare the CNN predictions with the dataset ground truth values, we set the CNN output to 6 values with 3 decimal places. We only experimented with a connector function who converting into characters the CNN output float but adding a spaces character between values and a separator character between the context and the target during training.

To optimize the first model, we adjusted the learning rate and dataset size. A dataset size of 200,000 samples for the CNN and 1,000,000 samples for the Transformer struck a balance between training time and performance. Additionally, we determined that a learning rate of  $1e-5$  was most efficient for the CNN, while a learning rate of  $3e-4$  yielded optimal results for the Transformer.

For the composite model hyper-parameters, we experimented two methods to connect the CNN outputs into a sequence of tokens: the first one was the same we use for chained models and another prioritizing important information such as the initial digits of each CNN float value. However, the latter approach introduced more issues than solutions, so we kept the former.

We also optimized the learning rate of the composite model and found that a value around  $1e-4$  yielded good results. When it came to the number of decimal places in the CNN outputs, experimenting with 3, 6, and 8 decimals revealed that 6 decimals provided the best performance. Three decimals did not provide sufficient relevant information to the Transformer, while 8 decimals overwhelmed it. Six decimal places struck the right balance.

Regarding the number of values provided by the CNN, increasing this count improved performance, and we settled on 76 values as the output. However, further increases in this value resulted in longer training times and higher resource requirements.

Finally, we explored the impact of dataset size on model optimization. Increasing the dataset size led to better results, but we were limited by time and available resources. Nevertheless, we were able to conduct experiments with 1,000,000 samples.

## 4 Results and discussion

Models	RMSE	Training time	Training labels
Chained models	$0.055 \pm 0.120$	4 hours	T*N
Composite model	$0.056 \pm 0.113$	2 days	N
CNN (baseline)	$0.061 \pm 0.055$	2.6 hours	N

Table 1: Results

In this paper we compared 3 models: the chained models, the composite model and the CNN.

The chained models implied to split the task into sub-tasks and required additional labeling effort to train each model. The training lasted 2 hours for each model to accurately predict the position of the square closest to the triangle with an accuracy of  $0.055 \pm 0.120$ .

The composite model allowed to avoid a large part of the labeling work by delegating the sub-task into a single so-called composite model. The training lasted 2 days and accurately predicted the position of the square closest to the triangle with an accuracy of  $0.056 \pm 0.113$ .

As a baseline, we also trained the CNN to extract only 2 values from the image, being the x y coordinates of the square closest to the triangle. We achieved an accuracy of  $0.061 \pm 0.055$  and a learning time of 2.6 hours.

Compared to the baseline, the two models studied are equivalent regarding the average RMSE values. However, they are slightly less efficient regarding the standard deviation. Keeping in mind our goal to address more complex use cases in the future (e.g, multi-modal models with text and images), both transformer-based models represent a suitable solution.

All these results are presented in the Table 1 and Figure 4 and we also quantify the amount of labeling work needed for each model, here T represents the number of sub-tasks that must be performed by our models and N, the amount of work to label a sub-task. The quantification of labeling work holds significant importance for us, particularly because our future work will involve larger datasets that are more complex and challenging to label compared to the ones discussed in this article. Therefore, our objective is to minimize this labeling work to optimize our time and enhance overall efficiency.

The results obtained from both the chained models and the composite model demonstrate that it is indeed possible to avoid the labeling work associated with training multitasking chained models, albeit at a notable additional cost. This finding is crucial for our future work as it enables us to concentrate more on fine-tuning the model itself, rather than solely focusing on the training datasets. By leveraging this information, we can allocate more resources towards optimizing and tuning the model for improved performance.

## 5 Conclusion

In this article, we compared chained models and composite model architectures. We observed that our chained models require more labeling effort than the composite model. Indeed, the composite model significantly reduces the labeling time by dividing it by the number of sub-tasks we aim to accomplish. Nevertheless, the composite model requires twelve times more training time to match the efficiency of the chained models.

Both architectures offer fair performances regarding our use case with only  $0.055 \pm 0.120$  and  $0.056 \pm 0.113$  RMSEs between the predicted and ground truth values. When sub-tasks can be explicitly defined and when training labels for each sub-task are available, the chained models architecture is preferred as it requires less hyper-parameters, saves training time and therefore power energy. Otherwise, the composite model architecture is more suitable when sub-tasks lack clear delineation and can't be clearly defined; additionally, it also offers the perspective to address future use cases requiring multi-modal inputs and various tasks.

Although further work is needed on other use cases to confirm our conclusion, these results prove essential for our future work, integrating multi-modal inputs and more elaborate images.

## References

- [1] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention, 2023.
- [2] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [3] OpenAI. Gpt-4 technical report, 2023.
- [4] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- [5] Abdellatif Bekkar, Badr Hssina, Samira Douzi, and Khadija Douzi. Air-pollution prediction in smart city, deep learning approach. *Journal of big Data*, 8(1):1–21, 2021.
- [6] Stefan Oehmcke, Lei Li, Jaime C Revenga, Thomas Nord-Larsen, Katerina Trepekli, Fabian Gieseke, and Christian Igel. Deep learning based 3d point cloud regression for estimating forest biomass. In *Proceedings of the 30th International Conference on Advances in Geographic Information Systems*, pages 1–4, 2022.

- [7] Ioannis D Apostolopoulos and Tzani A Mpesiana. Covid-19: automatic detection from x-ray images utilizing transfer learning with convolutional neural networks. *Physical and engineering sciences in medicine*, 43:635–640, 2020.
- [8] Chao-Jung Liu, Vladimir A Krylov, Paul Kane, Geraldine Kavanagh, and Rozenn Dahyot. Im2elevation: Building height estimation from single-view aerial imagery. *remote sensing*, 12(17):2719, 2020.
- [9] Alexander G Olenskyj, Brent S Sams, Zhenghao Fei, Vishal Singh, Pranav V Raja, Gail M Bornhorst, and J Mason Earles. End-to-end deep learning for directly estimating grape yield from ground-based imagery. *Computers and Electronics in Agriculture*, 198:107081, 2022.
- [10] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [12] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.
- [13] Aiden Nibali, Zhen He, Stuart Morgan, and Luke Prendergast. Numerical coordinate regression with convolutional neural networks, 2018.
- [14] Naveen Ram, Dima Kuzmin, Ellie Ka In Chio, Moustafa Farid Alzantot, Santiago Ontanon, Ambarish Jash, and Judith Yue Li. Multi-task end-to-end training improves conversational recommendation, 2023.
- [15] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning*, pages 5583–5594. PMLR, 2021.
- [16] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020.
- [17] Wanjie Lu, Chaozhen Lan, Chaoyang Niu, Wei Liu, Liang Lyu, Qunshan Shi, and Shiju Wang. A cnn-transformer hybrid model based on cswin transformer for uav image object detection. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2023.
- [18] Hussam Qassim, Abhishek Verma, and David Feinzimer. Compressed residual-vgg16 cnn model for big data places image recognition. In *2018 IEEE 8th annual computing and communication workshop and conference (CCWC)*, pages 169–175. IEEE, 2018.

## Appendix: distance calculation with GPT4.0

This section presents a dialogue exchange (prompts and responses) based on multimodal inputs (text and images) performed with the GPT-4.0 model (as of April 14, 2023). For better readability, the URL prefix of PNG images has been removed from the dialogue.

- *What does the image Full\_GPT\_architecture.png represent?*  
The image represents the architecture of a Full GPT (Generative Pre-trained Transformer) model. GPT is a type of transformer-based language model used for various natural language processing tasks, such as text generation, translation, and summarization. The image shows the different layers and components of the GPT model, including the input and output layers, self-attention mechanism, position-wise feed-forward networks, and layer normalization.
- *In image map\_0.png the point with coordinates  $x=0.000$  and  $y=0.000$  is at the bottom left of the image and the point with coordinates  $x=1.000$  and  $y=1.000$  is at the top right of the image. The triangle is positioned at  $x=0.622$  and  $y=0.439$ ; a first square is positioned at  $x=0.517$  and  $y=0.898$  and a second square is positioned at  $x=0.378$  and  $y=0.886$ . The closest square to the triangle is positioned at  $x=0.517$  and  $y=0.898$ . Is that clear for you?*  
Yes, that is clear. The triangle is at coordinates (0.622, 0.439), and there are two squares - one at (0.517, 0.898) and the other at (0.378, 0.886). The closest square to the triangle is the one at (0.517, 0.898).
- *In image map\_1.png the point with coordinates  $x=0.000$  and  $y=0.000$  is at the bottom left of the image and the point with coordinates  $x=1.000$  and  $y=1.000$  is at the top right of the image. The triangle is positioned at  $x=0.577$  and  $y=0.713$ ; a first square is positioned at  $x=0.778$  and  $y=0.263$  and a second square is positioned at  $x=0.583$  and  $y=0.137$ . The closest square to the triangle is positioned at  $x=0.778$  and  $y=0.263$ . Is that*

clear for you?

Yes, that is clear. In the image, the triangle is at coordinates (0.577, 0.713), and there are two squares - one at (0.778, 0.263) and the other at (0.583, 0.137). The closest square to the triangle is the one at (0.778, 0.263).

- *In image map\_2.png the point with coordinates  $x=0.000$  and  $y=0.000$  is at the bottom left of the image and the point with coordinates  $x=1.000$  and  $y=1.000$  is at the top right of the image. The triangle is positioned at  $x=0.127$  and  $y=0.626$ ; a first square is positioned at  $x=0.196$  and  $y=0.341$  and a second square is positioned at  $x=0.152$  and  $y=0.243$ . The closest square to the triangle is positioned at  $x=0.196$  and  $y=0.341$ . Is that clear for you?*

Yes, that is clear. In the image, the triangle is at coordinates (0.127, 0.626), and there are two squares - one at (0.196, 0.341) and the other at (0.152, 0.243). The closest square to the triangle is the one at (0.196, 0.341).

- *In image map\_3.png the point with coordinates  $x=0.000$  and  $y=0.000$  is at the bottom left of the image and the point with coordinates  $x=1.000$  and  $y=1.000$  is at the top right of the image. What is the position (x, y) of the square closest to the triangle?*

In the image map\_3.png, the triangle is positioned at (0.372, 0.497). There are two squares - one at (0.430, 0.542) and the other at (0.283, 0.464). The closest square to the triangle is the one at (0.283, 0.464).

- *No, in the image map\_3.png the triangle is positioned at  $x=0.597$  and  $y=0.730$ ; a first square is positioned at  $x=0.693$  and  $y=0.873$  and a second square is positioned at  $x=0.323$  and  $y=0.142$ . The closest square to the triangle is positioned at  $x=0.693$  and  $y=0.873$ .*

Apologies for the confusion. In the image map\_3.png, the triangle is positioned at (0.597, 0.730). There are two squares - one at (0.693, 0.873) and the other at (0.323, 0.142). The closest square to the triangle is the one at (0.693, 0.873).

- *In image map\_4.png the point with coordinates  $x=0.000$  and  $y=0.000$  is at the bottom left of the image and the point with coordinates  $x=1.000$  and  $y=1.000$  is at the top right of the image. What is the position (x, y) of the square closest to the triangle?*

In the image map\_4.png, the triangle is positioned at (0.497, 0.372). There are two squares - one at (0.542, 0.430) and the other at (0.464, 0.283). The closest square to the triangle is the one at (0.464, 0.283).

- *Yes, you had almost the right answer. In image map\_4.png, the triangle is positioned at  $x=0.138$  and  $y=0.135$ ; a first square is positioned at  $x=0.564$  and  $y=0.311$  and a second square is positioned at  $x=0.597$  and  $y=0.211$ . The closest square to the triangle is positioned at  $x=0.564$  and  $y=0.311$ .*

Apologies for the confusion. In the image map\_4.png, the triangle is positioned at (0.138, 0.135). There are two squares - one at (0.564, 0.311) and the other at (0.597, 0.211). The closest square to the triangle is the one at (0.564, 0.311).

- *In image map\_5.png the point with coordinates  $x=0.000$  and  $y=0.000$  is at the bottom left of the image and the point with coordinates  $x=1.000$  and  $y=1.000$  is at the top right of the image. What is the position (x, y) of the square closest to the triangle*

In the image map\_5.png, the triangle is positioned at (0.497, 0.372). There are two squares - one at (0.542, 0.430) and the other at (0.464, 0.283). The closest square to the triangle is the one at (0.464, 0.283).

- *No, in image map\_5.png the square closest to the triangle is positioned at  $x=0.604$  and  $y=0.148$ .*

Apologies for the confusion. In the image map\_5.png, the triangle is positioned at (0.604, 0.148). The closest square to the triangle is also positioned at (0.604, 0.148).

To replay the prompts, use the following URL prefix: "<https://raw.githubusercontent.com/marxav/aloo/main/>".

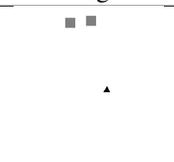
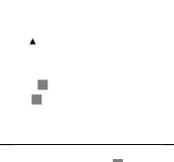
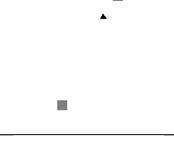
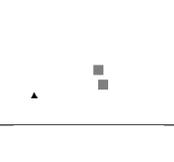
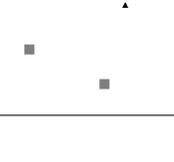
Image	image_name	triangle	squares[0]	squares[1]	target
	map_0.png	(0.622, 0.439)	(0.517, 0.898)	(0.378, 0.886)	(0.517, 0.898)
	map_1.png	(0.577, 0.713)	(0.778, 0.263)	(0.583, 0.137)	(0.778, 0.263)
	map_2.png	(0.127, 0.626)	(0.196, 0.341)	(0.152, 0.243)	(0.196, 0.341)
	map_3.png	(0.693, 0.873)	(0.323, 0.142)	(0.597, 0.730)	(0.693, 0.873)
	map_4.png	(0.564, 0.311)	(0.597, 0.211)	(0.138, 0.135)	(0.564, 0.311)
	map_5.png	(0.103, 0.382)	(0.604, 0.148)	(0.743, 0.673)	(0.604, 0.148)

Table 2: Images used to experiment with GPT4.0