

D.I.A.G.R.A.M.: Development of Image Analysis for Graph Recognition And Modeling

Filippo Garagnani, Saverio Napolitano, Nicola Ricciardi

'Computer Vision and Cognitive System' course

Università di Modena e Reggio Emilia

Abstract—This report describes D.I.A.G.R.A.M., a system for diagram recognition and generation. We present its architecture, the logic behind deterministic algorithms and the training of its deep learning components. The project aims to transform visual diagram input into structured representations.

I. INTRODUCTION

Diagrams are crucial in education, documentation, and many other fields. Developing a system that automatically understands and generates diagrams may reveal helpful for having high-quality, easily editable representations. This, however, poses challenges involving classification, shape detection, structure interpretation and symbolic representation. Our project proposes a modular architecture to tackle these tasks.

II. PREVIOUS WORKS

- [1]
- [2]
- [3]

III. SYSTEM ARCHITECTURE

Figure 1 illustrates the full pipeline, composed of four main modules: the Classifier, Extractor, Transducer, and Compiler. Each module is designed to process and transform the handwritten diagram image progressively toward a structured output.

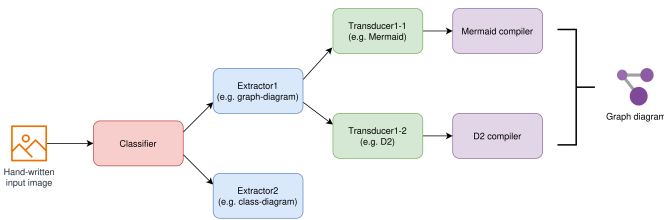


Fig. 1: Overview of the D.I.A.G.R.A.M. system architecture.

The Classifier is able to recognize the category of the hand-written diagram that the user submitted as input. This is necessary in order for the system to know to which Extractor pass the data. This module is able to apply object detection and semantic recognition to represent the diagram in a unified way. After that, the representation is sent to the Transducer, which converts the data in a Markup language of choice. The Markup language content is therefore sent to the Compiler, thus generating the high-quality and editable diagram in .png format.

A. Classifier Module

The Classifier is the first component of the pipeline. It is able to tell to which category a given image of an handwritten diagram belongs to. This is necessary in order to later know the Extractors that can process the input image.

B. Extractor Module

The Extractor is the key component of the D.I.A.G.R.A.M. system, able to transform the raw diagram image into a structured, category-agnostic representation. The Extractor recognizes the diagram's components, such as shapes, lines, and text, through an object detection network, and organizes them into a unified format. This representation serves as the foundation for the subsequent Transducer module, which converts the structured data into a domain-specific markup language.

C. Transducer Module

The Transducer is responsible for converting the unified, agnostic representation of a diagram into a domain-specific markup language. This transformation enables the subsequent Compiler module to generate high-quality visual outputs.

D. Compiler Module

The Compiler module is the final stage of the D.I.A.G.R.A.M. pipeline. Its primary role is to take the structured representation of the diagram, expressed in a markup language (e.g., Mermaid.js), and generate a high-quality diagram in a visual format such as PNG.

IV. FLOWCHART DIAGRAM RECOGNITION

A. Internal Representation

In order to pass data through the components of the pipeline in a unified way, it was necessary to decide upon an internal representation of the flowchart diagram. It consists of **elements** and **relations**. The latter are linked to at least one element. Every element has a *category* and possibly some *inner* and *outer text*. Every relation has a *category* too, either a *source* or *target element* or both, and possibly *text* related to it.

This internal representation is generated from the image provided by the user by the Extractor module, while it is parsed and converted to a markup language of choice by the Transducer.

B. Classifier Module

1) *Preprocessing*: Due to the fact that various kinds of diagrams' images could be submitted to the system, it was decided to streamline a unified preprocessing pipeline in order to reduce the differences that our classifier model had to handle. A few problems were analyzed and the final preprocessing pipeline for the classifier module is composed of the following operators, in order:

- **Gray Scaling**: In order to remove colors from images, which was deemed an useless information for classifying diagrams.
- **Otsu Thresholding**: To binarize the image and reducing the information to handle.
- **Median Filtering**: Some images in our dataset - and in the use-case - were taken on checkered notebooks. The application of a median filter, with kernel size equal to three, totally removed these.
- **Perspective Correction**: The diagrams' images, sometimes, didn't occupy the full photo frame. Finding four keypoints (if possible) and stretching the image helped to fully employ the frame.
- **Padding**: This was necessary in order to keep the same net for every image. Changes in resolution would have implied changes in the number of paramters of the MLP layers.

Figure 2 shows the step-by-step process of preparing an input image for the network.

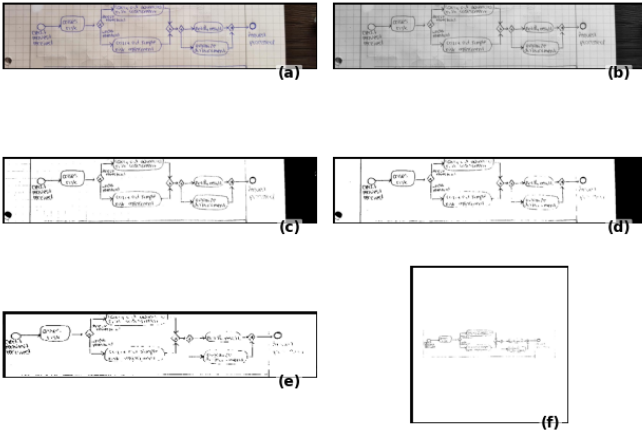


Fig. 2: The operations applied by the preprocessing module. (a): The raw image, still colored. (b): The image in gray scale. (c): The binarized image, thanks to Otsu. (d): The median filter result. Some previous thin lines are now no more. (e): The image with a perspective correction. (f): The image in a standard size, ready to enter the network.

2) *Model*: Since it was assumed that classifying two types of diagrams - with a third 'other' class - was a quite simple task, a small network was deemed enough. The model consists of three convolutional blocks - each containing a convolutional layer, followed by a batch normalization and a

max-pooling layer - then followed by three linear layers (for a more detailed overview, see Appendix C). The output vector of the MLP contains three features - one for each class that this specific classifier needs to address. In order to determine the prediction, the softmax function is then applied. As later seen in section VI-A, a weighted Cross-Entropy loss is then used to train the model.

C. Extractor Module

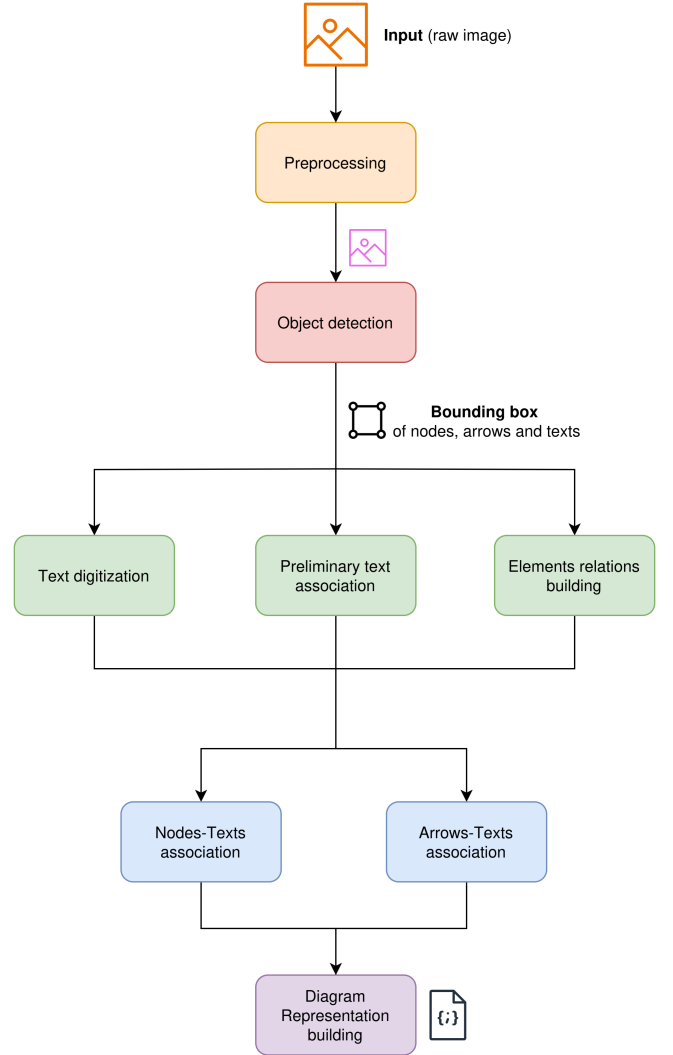


Fig. 3: The complete pipeline of the extractor module.

1) *Preprocessing*: After some brief experiments, we found out that the object detection task through bounding boxes worked very well without any preprocessing technique. Moreover, the use of certain filters actually *downgraded* the performances - e.g. the use of aggressive median filters reduced the capacity of the OCR model for text extraction.

The only preprocessing applied is a gray-scale transformation. Other transformations are not deemed necessary - the model employed automatically resize the

input image, so it's not needed to do it beforehand.

2) *Bounding Box Detection*: TODO: FasterRCNN - backbone ResNet 50FPN [4].

After having determined all the bounding boxes of the diagram image, in order to correctly create the diagram internal representation, there's the need to link together text, elements and relations.

3) *Arrow Orientation*: To later correctly assign arrows to their source and target elements, the need to determine where the tip of an arrow lies. This task involves a few problems, among which determining exactly the arrow shape and its end - because sometimes other pen signs may be found in the bounding box of the arrow itself. After a few trials (and errors; for further insight, refer to section VI-C), the approach yielding the best results was through the use of the object detection network (see section IV-C2) to recognize both arrow heads and tails.

4) *Text Digitization*: Since the task of converting an image of text into the relative string has already been solved in various ways, we settled upon the Microsoft *Trocr Handwritten* Transformer model; more precisely, the *base* version of it [5]. Every bounding box containing text that has been found is later passed to the model, and the resulting string generated by the model is linked to it.

5) *Relationships Building*: Once the text has been digitized, it's necessary to effectively compute which elements are connected by which arrows. This is done by finding, for each arrow, the elements nearer to the head (which we will assume to be the 'target' of the relation) and to the tail (the 'source' of the relationship exemplified by the arrow). These associations are nullified if the distance is above a certain adjustable threshold. This last step is done under the hypothesis that an arrow may come from nothing - or may point to nothing, having actually only either the source or the target.

6) *Text Association*: It can be assumed that every piece of text is linked to one and only one element or relation. So, every bounding box containing text is linked to the nearest element - it being either a relation or an element. The associations that require too long a distance are eliminated further in the algorithm, and the relevant text is discarded.

7) *Element-Text Association*: After that a text box has been assigned to an element, there's the need to understand if the text is either inside or outside the element and whether the text is actually referring to the given element or not. In order to do this, the overlap between the two bounding boxes is computed; if the overlap is over a certain threshold, heuristically decided, the association is kept. Otherwise, the distance between the two bounding boxes is computed - as the distance between their central point. If the distance is

over a certain threshold, the association is kept; otherwise, it's discarded.

8) *Relation-Text Association*: After a text box is assigned to a relation, some computation is needed to fully understand how the text is associated to the relation itself. The relation's bounding box is then split into three different sections. The overlap between the text bounding box and each of the three sections boxes is computed - then, the section with the highest overlap value is then assigned to the text under consideration.

V. DATASET

TODO: Elencare i dataset usati, in che modo erano usati originariamente e come lo usiamo noi (quali dati per quali componenti)

VI. EXPERIMENTS

A. Classifier

The first architecture explored for the Classifier module was a Convolutional Neural Network (CNN) following AlexNet's architecture, but smaller in size. It had two Convolutional layers, each followed by a batch normalization and then by a max pooling layers. The first layer had only 4 filters, while the second one had 8. The output of the second layer was then flattened and passed through a fully connected layer which downsized the features to 128, followed by a ReLU activation function. Other two fully connected layer brought the number of output features down to the number of classes. Then, a softmax activation function was applied to the output layer, in order to obtain the probabilities of each class. The model was trained on our own dataset, with only three possible output classes: 'flowchart', 'graph' and 'other'; it was trained for 10 epochs, with an AdamW optimizer (with a learning rate of $1e-3$) and a standard CrossEntropy Loss function.

We first tried such a small network because the features that had to be recognized were very simple, down to the basic shapes of the diagram elements. However, the results showed a non-decreasing value of the loss function of 0.9 - better than $-\ln(3) \simeq 1.1$, which is random guessing, but still not satisfactory.

Judging this to be a problem in architecture size, we then tried adding another CNN layer with 16 filters; this one too followed by a batch normalization and a max pooling layer. The rest of the net was untouched; we also decreased the learning rate to $1e-4$, in order to avoid overshooting the minimum of the loss function. The results, however, didn't show any kind of improvement, with the loss function still being around 0.9.

While analyzing the predictions, we came up with what was actually the problem: the dataset was heavily unbalanced, with the 'other' class being the most represented one, shadowing the

other two classes. In order to solve this problem, we decided to use a weighted CrossEntropy Loss function, with the weights being inversely proportional to the number of samples in each class.

Algorithm 1 Cross-Entropy weight computation

Require: The classes' frequencies f_1, f_2, \dots, f_c

```

1:  $S \leftarrow \sum_{i=1}^c f_i$ 
2:  $f_i \leftarrow (f_i/S)^{-1}$ 
3:  $S \leftarrow \sum_{i=1}^c f_i$ 
4:  $f_i \leftarrow f_i/S$ 
5: return  $f_1, f_2, \dots, f_c$ 

```

This final approach showed a significant improvement in the results, with the loss function decreasing to around 0.6 after 10 epochs of training, with a global accuracy score of 97% in the validation set.

B. Object Detection

TODO: Talk about the net.

The model employed was tested in two different ways: as zero-shot, after a 10-epoch training and after a 100-epoch training. The best-behaving one was, of course, the latter. In Table I and Table II we show the Mean Average Precision and Recall of the model used.

TABLE I: Mean Average Precision (*mAP*) Metrics

Metric	Value	Description
mAP	0.8754	Overall mean Average Precision
mAP@0.5	0.9878	mAP at IoU threshold 0.5
mAP@0.75	0.9183	mAP at IoU threshold 0.75
mAP (Small)	0.6444	mAP for small objects
mAP (Medium)	0.8054	mAP for medium objects
mAP (Large)	0.8975	mAP for large objects

TABLE II: Mean Average Recall (*mAR*) Metrics

Metric	Value	Description
mAR@1	0.3436	mAR with max 1 detection per image
mAR@10	0.8877	mAR with max 10 detections per image
mAR@100	0.9128	mAR with max 100 detections per image
mAR (Small)	0.6631	mAR for small objects
mAR (Medium)	0.8670	mAR for medium objects
mAR (Large)	0.9141	mAR for large objects

Furthermore, we also think that the following metrics show a great precision:

- **Average X-axis error:** 6.84 pixels
- **Average Y-axis error:** 4.84 pixels

C. Arrow Orientation

As we've seen in the previous section, the need to find the exact position of an arrow's head arises (see **Appendix A** and **Appendix B** for a full overview of the first approaches).

The solution, as stated before, was to finetune the already working RCNN in order for it to also extract the bounding boxes of both tails and heads of the arrow.

D. Text Digitization

In order to choose an appropriate model, some tests were done over our specific dataset employing some different, pre-made models. The average results over the dataset are shown in Table III.

TABLE III: Text Distance Metrics

Model	Hamming	Cosine	Euclidean
microsoft-trocr-small-printed	1.392	0.150	7133
microsoft-trocr-small-handwritten	1.346	0.100	3947
microsoft-trocr-base-handwritten	1.549	0.057	3003

Some examples of the different models applied to our flowchart diagram dataset can be seen in Figure 4 and in Figure 5.

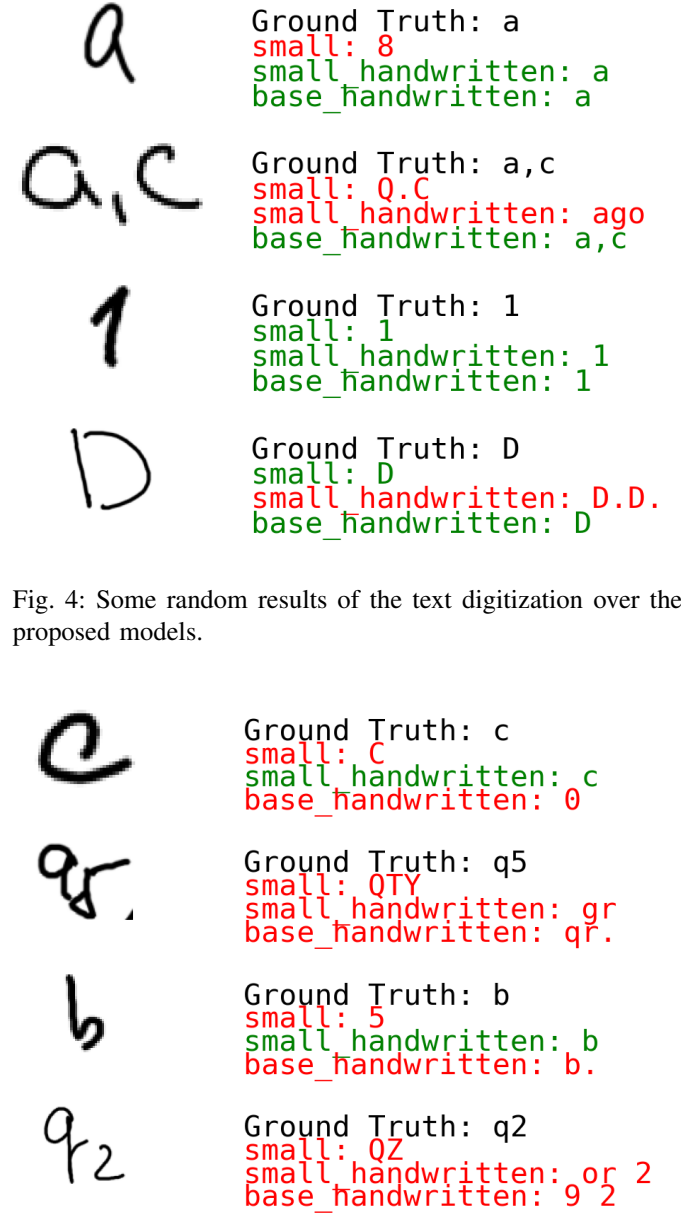


Fig. 5: Some results in which the *trocr-base-handwritten* fails.

VII. DISCUSSION

TODO

VIII. CONCLUSION AND FUTURE WORK

TODO

REFERENCES

- [1] B. Schäfer, M. Keuper, and H. Stuckenschmidt, "Arrow r-cnn for handwritten diagram recognition," *International Journal on Document Analysis and Recognition (IJDAR)*, vol. 24, pp. 3–17, 2021. [Online]. Available: <https://doi.org/10.1007/s10032-020-00361-1>
- [2] E. Lank, J. Thorley, and S. Chen, "An interactive system for recognizing hand drawn uml diagrams." 01 2000, p. 7. [Online]. Available: <https://doi.org/10.1145/782034.782041>
- [3] I. Saito, H. Yoshida, and K. Sakaguchi, "Sketch2diagram: Generating vector diagrams from hand-drawn sketches," in *The Thirteenth International Conference on Learning Representations*, 2025. [Online]. Available: <https://openreview.net/forum?id=KvaDHPPhir>
- [4] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, Eds., vol. 28. Curran Associates, Inc., 2015. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2015/file/14bfa6bb14875e45bba028a21ed38046-Paper.pdf
- [5] M. Li, T. Lv, J. Chen, L. Cui, Y. Lu, D. Florencio, C. Zhang, Z. Li, and F. Wei, "Trocrr: Transformer-based optical character recognition with pre-trained models," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 11. AAAI Press, 2023, pp. 13 094–13 102. [Online]. Available: <https://doi.org/10.1609/aaai.v37i11.26538>
- [6] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004. [Online]. Available: <https://doi.org/10.1023/B:VISI.0000029664.99615.94>
- [7] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 888–905, 2000. [Online]. Available: <https://doi.org/10.1109/34.868688>
- [8] D. Deng, "DbSCAN clustering algorithm based on density," in *2020 7th International Forum on Electrical Engineering and Automation (IFEEA)*, 2020, pp. 949–953. [Online]. Available: <https://doi.org/10.1109/IFEEA51475.2020.00199>
- [9] N. Siddique, S. Paheding, C. P. Elkin, and V. Devabhaktuni, "U-net and its variants for medical image segmentation: A review of theory and applications," *IEEE Access*, vol. 9, pp. 82 031–82 057, 2021. [Online]. Available: <https://doi.org/10.1109/ACCESS.2021.3086020>

APPENDIX A DOUBLE CLUSTERING

The first algorithmic approach proposed to solve this task - called *Double Clustering* - was as follows:

- 1) Apply a clustering algorithm to the bounding box.
- 2) Suppose that the biggest cluster belongs to the arrow; the other clusters are thought to contain useless elements for this task, and can be discarded.
- 3) Find the keypoints of the arrow using SIFT [6].
- 4) Find the two keypoints inside the arrow with maximum distance. Assume that these two points are closely related to the two ends of the arrow.
- 5) Run a template matching procedure inside the two clusters to find the head of the arrow.

For Step (1), we tried to employ both Spectral Clustering [7] (see Figure 6) and DBScan [8] (see Figure 7).



Fig. 6: The result of Spectral Clustering to one instance of an arrow bounding box in our dataset.



Fig. 7: The result of DBScan to one instance of an arrow bounding box in our dataset.

For Step (3), we used SIFT [6] in order to find keypoints and some descriptors for them - necessary for the later template matching. During Step (4), we were able to correctly deduce the overall direction of the arrow (see Figure 8).

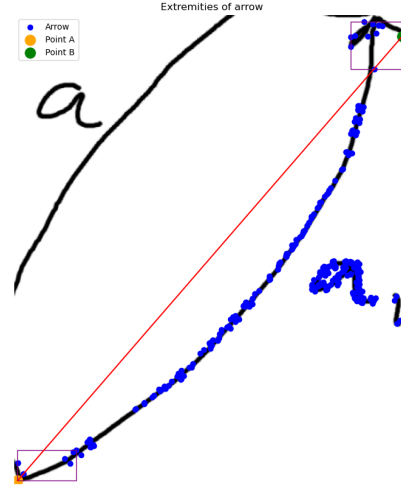


Fig. 8: The estimated direction of the arrow, finding the points that maximize the distance inside the cluster.

For Step (5), we had previously hand-picked some templates (see Figure 9).



Fig. 9: An example of hand-picked template with SIFT descriptors.

However, this approach didn't work as expected (see Figure 10). In particular, it never worked around self arrows: the maximized distance inside one of the two clusters is never between the tail and the head of the arrow.

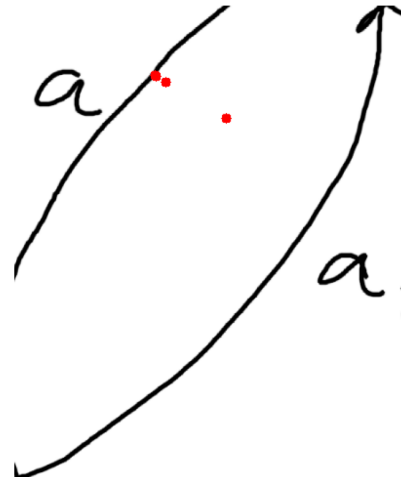


Fig. 10: The overall result of the algorithm.

APPENDIX B ARROWNET

This second approach for determining an arrow's orientation was inspired by UNet [9]. The core idea was to segment an arrow's bounding box, and to use a specific CNN to try and determine the arrow's head point and the arrow's tail point. This downsampling-upsampling network is called **ArrowNet** and its specifics can be found in Table IV.

TABLE IV: ArrowNet Architecture Details

Stage	Layer	Out Ch	Operations
Encoder	Encoder1	32	Double Conv Block
	MaxPool	32	2×2 MaxPooling
	Encoder2	64	Double Conv Block
Bottleneck	MaxPool+Bottleneck	128	MaxPool+Double Conv
Decoder	Upsampling2	64	Transpose Conv Block
	Skip + Decoder2	64	Concat + Double Conv
	Upsampling1	32	Transpose Conv Block
	Skip + Decoder1	32	Concat + Double Conv
Output	Output Conv	2	1×1 Convolution

The entire architecture is composed of three main building blocks:

- 1) **Double Convolution Block:** Each double convolution block applies two consecutive 3×3 convolutions with padding 1, followed by batch normalization, and a ReLU activation after each convolution. This should preserve spatial dimensions while increasing feature depth.
- 2) **Upsampling Block:** The upsampling blocks use 2×2 transposed convolutions with stride 2, followed by batch normalization and ReLU activation.
- 3) **Skip Connections:** The architecture employs skip connections that concatenate feature maps from the encoder path with corresponding decoder features.

The two output float channels of ArrowNet should have been used to predict, in a given position: 1 if the head's point is in a certain position, 0 else; 1 if the tail's point is in a certain position, 0 else - for the first and second channel respectively.

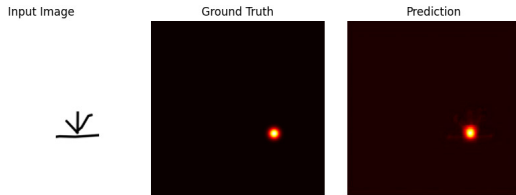


Fig. 11: Some results for arrow's head reconstruction.

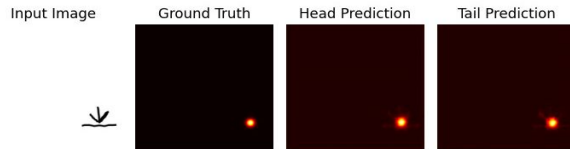


Fig. 12: The overall result of an ArrowNet input.

As seen in Figure 12, due to the low number of samples in the dataset, the network is unable to correctly distinguish between an arrow's head and tail.

APPENDIX C THE CLASSIFICATION NETWORK

The **Classification Network** is a lightweight convolutional neural network, suited for processing grayscale images. The architecture is composed of convolutional feature extraction layers followed by fully connected classification layers. Some specific about it can be found in Table V.

TABLE V: Classification Network Details

Layer Type	Out Channels	Kernel Size	Pool Kernel Size
Conv	8	2	3
Conv	16	2	3
Conv	24	3	2
Layer Type	Out Features	Dropout Rate	Activation Function
Linear	2048	0.2	ReLU
Linear	128	0.2	ReLU
Linear	3	0.2	ReLU