### Transforming IceCube

*Energy and Angular Regression of IceCube Data using a NN, GNN, and Transformer Architecture* 

Laurits Møberg, Franciszek Nawrocki & David Johansen 11.06.2025

(All parts contributed equally)



UNIVERSITY OF COPENHAGEN

### Outline

- > Data
- ➤ Tools
- Neural Network
- ≻ GNN
- > Transformer
  - > Energy
  - Direction

### Conclusion

Source:



Simulated Data – No outliers/NaNs Energy range: 10<sup>4</sup> - 10<sup>6</sup> GeV

**PMT-fication** 

124/dom Multiple pulses -> 32/dom Variables

Pulses -> List

Events: ~200k per file

1	🕕 [Receipt]22011_1.json	💶 truth_1.parquet
2 2	🕕 [Receipt]22011_2.json	💶 truth_2.parquet
3	🕕 [Receipt]22011_3.json	💶 truth_3.parquet
1 4	[] [Receipt]22011_4.json	🚺 truth_4.parquet
5	🕕 [Receipt]22011_5.json	💶 truth_5.parquet
<mark> </mark> 6	🕕 [Receipt]22011_6.json	💶 truth_6.parquet
7	🕕 [Receipt]22011_7.json	💶 truth_7.parquet
8 📕	🕕 [Receipt]22011_8.json	💶 truth_8.parquet
9	🕕 [Receipt]22011_9.json	💶 truth_9.parquet
10	🕕 [Receipt]22011_10.json	💶 truth_10.parquet

### Data

• Dataloader

class PMTfiedDatasetPyArrow(Dataset): def \_\_init\_\_( self, truth\_paths, selection=None, transform=feature\_preprocessing, ):

def

• Feature preprocessing

make_dataloader_PMT <mark>(</mark>	
root_dir,	
dataset_id,	
training_parts,	
validation_parts,	
<pre>batch_size=config['training_params']['batch_size'</pre>	],
<pre>num_workers=config['input_data']['num_workers'],</pre>	

<pre>if col_name in ['dom_x', 'dom_y', 'dom_z', 'dom_x_rel', 'dom_y_rel', 'dom_z_rel']:     value = value / 500</pre>
elif col_name in ['rde']:
value = (value - 1.25) / 0.25
elif col_name in ['pmt_area']:
value = value / 0.05
elif col_name in ['q1', 'q2', 'q3', 'q4', 'q5', 'Q25', 'Q75', 'Qtotal']:
mask = value > 0
<pre>value[mask] = np.log10(value[mask])</pre>
elif col_name in ['t1', 't2', 't3','t4', 't5']:
mask = value $> 0$
value[mask] = (value[mask] - 1.0e04) / 3.0e04
elif col_name in ['T10', 'T50', 'sigmaT']:
mask = value $> 0$
value[mask] = value[mask] / 1.0e04

## Data







### Tools

Data Tools NN GNN Transformer Conclusion

# **Need** GPU for efficient training Solution:



### Workflow:

Code locally -> Push to Github

-> Clone and run on Colab

[ ] !git clone <u>https://github.com/Moebergs/AppML\_FinProject</u> %cd AppML\_FinProject/training\_and\_inference/

] !pip install -r requirements.txt -q

[ ] from google.colab import drive drive.mount('/content/drive')

#### import zipfile

zip\_ref = zipfile.ZipFile("/content/drive/MyDrive/22011\_1T4.zip", 'r')
zip\_ref.extractall("/content/dataset")
zip\_ref.close()

#### [ ] !python -m train

Tools

Even with Colab, training the models takes a lot of **time**.

Total time for testing and training:

~ 250 Hours

Need way to monitor training Solution:

Weights & Biases (wandb.ai)

## Weights & Biases



### **Neural Network**





### It is learning...



### GNN



#### Transforming IceCube

### GNN



### Transformer

### Attention Is All You Need



Output prediction

### Transformer

Transformer encoder model made by Luc Voorend (https://github.com/lucvoorend/IceCubeTransformer)



Data Tools NN GNN Transformer Conclusion

From *Angular reconstruction of high energy neutrinos using machine learning* by L. Voorend, 2025, University of Copenhagen.

### Transformer – Energy Regression

Trained in log domain Model is learning, but it is exhibiting a strong bias towards the mean





### Transformer – Energy Regression

Mean predicted points clearly correlated with zenith and number of domes.

Quick test: Add number of domes as input to final linear layer. Little to no improvement

0.5

0

-0.5

-<del>1</del> -0.5

GNN Transformer

Data

Tools

NN

Conclusion

Idea: Train model for two hemispheres separately



### Transformer – Energy Regression











### Transformer – Energy Regression – Further work







Aggregation

Mean, Sum, Min and Max instead of just Mean

**Final regression** 

Replace final linear layer with more complex MLP.

Implemented both in model with no increase in performance (See Appendix). However, no optimization was done. Data Tools NN GNN Transformer Conclusion

### Transformer – Direction

x, y, z →θ, φ

- Von Mises Fisher Loss  $\log \operatorname{prob}(x_{\operatorname{truth}} | x_{\operatorname{predicted}}, \kappa)$
- к Concentration/certainty parameter
- Calculated using 3vectors with norm 1



### Transformer – Direction

### **Azimuthal Angle**

### Zenith Angle



Periodic Domains – Transformation of predicted values to account for physical interpretation

## Data Tools GNN Transformer Conclusion



### Azimuthal Angle

### Zenith Angle



### Curiosity: Attention Weights - A window into the black box

### Conclusion

- Positive experience with 'the Troels model'
- Multiple models to compare
- Room for optimization with energy
- Significant difference in prediction power for different targets

## Thank you for listening!

### Special thanks to Luc for code and assistance (All parts contributed equally)

Links to our github pages: <u>https://github.com/DavidJohansen/IceCubeTransformer</u> <u>https://ghp\_mHATxjWQiUDwAS1AkTK4FdKsUYAt0B2eQB9f@github.com/mister-fran/FinalProjectFran</u> https://github.com/Moebergs/AppML\_FinProject

## Appendix

### Appendix A – Data

All data is generated through simulations. This means that it contains no NaNs, outliers or any other difficulties.

This does however give certain limitations to the applicability of the models. Domain-shifts will be expected upon transition from simulated to real-world data.

A few parameters are only simulated in a narrow parameter space. The parameter bad\_dom\_status is only simulated with a value of 1, meaning that any deviation from this value can give rise to unexpected predictions. This is also the case for a few other parameters.

### Appendix A – Data

The simulated data does not have a spherical distribution, due to the neutrinos' interaction with the Earth and the atmosphere. We do, however, have uniform azimuth distribution, giving rotational symmetry.



### 3d scatter plot of spherical distribution of data points



### Appendix A – Data (feature correlation and distribution)

200 200 					1 E	•																100 E						
											4 <sup>1</sup>													R			<b>8</b> . 8	
2	V. James		¥ 4	N. K.	. 9	:															a fait				19 (Aug			
1	<b>Ultim</b>	În alt		a and a start	-			for a second	bourse .	Sec. Sec.	· · · ·		the -		in						Sec.		and the second s	12	1985	Second 1	S. 64	and the second second
	<b>MARKET</b>				-				and a second		a inte	<u>κ</u> γ :	Sec. 1.	and the second	and the second					1	and the second s	1	i si		i an	and the second s	B. An	
1. 1984	William		- stocks-			1		ters a second	and the second	Sec. Sec. 1	and a second	6.> :	Mary 1.	Section 1.	and the second			1			Mar Carl	No.	liger.	hier	Mercher	and the state of the	Sec. 14.15	and the second
					-			-	-							1 1		1					!			<u> </u>		·
																		• • •										
	1.17.								1.1		12	· · · · · · · · · · · · · · · · · · ·			1					1				1	5			
i tadhi tadin bi alia	madadatatakinika	Advantation of the							4687 <u>, 1997</u> ,	aller		<b>86.7</b>			Hield - All						SPAteman and a							H.L. Hilds. Sam.
	anhoithlatan	Andreas Statistics	an combattle Bittlesin por	a naplati Planina					1	UNDER .		1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		2000 1	ANA									Pitchanna II	Bith			
Ladenirel C.C. Barr	adkambilishi	and the second second second	and a fill filling a	a naphiriti Dikanka	- paint and the major of							1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1			ABARCS V										1 1 1	Deserved on p to a served of the served of t	1000000	Article Ballinester
	attanidation a			La manatafahanan.							1	107 -			administration		!	· · · · · ·				388	Stanna .					anatellina
5	an atil at and	1 <u>Alberta Barboa</u> .		an an in the second	1.			36		alla -	1997 1997 1997		alle.9 2	and a set of	addets s						1. 1.		1. 1.		10000000000000000000000000000000000000			minad Billin
1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1	ana	and the second second	1							14 · · · ·	1000 - 10000 - 10000 - 1000 - 1000 - 1000 - 1000 - 1000 - 1000 - 1000 -		14		and Maria and Anna an							- Minana -	- Minunger -					A
internation in the	and diates	diam dedica	S S						1997 - 19		900 2010 - 11 2010 - 11		1917 - 11 2017 - 11		-atili-					l					1944	: ; ;		adarahita
	ditte of the distribution	Antonia della d							<u>MG 27</u>	<u>m</u>	<b>»</b> -	<b>1</b>	1978 ·				!	!	· · · · · · · · · · · · · · · · · · ·	<u>.</u>	Bhile	Mdn	Nda		Side			pititiki ili
			L.,					-			-	-		-				· · · · · ·									<u></u>	
							:	-	-						-			• • • • •							<u></u>		!	
						<u>.</u>		-	-	-				-	-	· · · ·		I										
-6.0 1.0 1.0								-										······	I		-							
-0.5 -1.0 -1.0 -1.0 -1.0 -1.0 -1.0 -1.0 -1.0		ALL ALL	Contraction of the second		3622				L.					-														
	APR AN	104-72-1		Same	and the second			La contra	Reduce of	<u>6.</u>	Read of the		Maria a c	India and a	Martin						THE .	AT 36	Mesie Mar	A W	the side	and the second second	and the second	Mathematica
	abire distantes		and the second					Ĩ.													22			130	Sec. 20.	Televis Contractor	0	all faithful the state
aris Statuta din S	anna shillided	Service and	State States	Selection of the second	- Art									line .							and a	200			A.S.	Tellas (S. M		(c)(2) Millionationkers
	and distant		2.45 <sup>2</sup>	1-1-1-2 	and the second				la contraction de la contracti	i file			-	Inc.									*	2		San year	\$	10 Charling
000	1.1.1.1.1.1	1.1							1											1			No.				1.	â :
- AND	olai nidorit	Real Maria							1	<b>L</b>	<b>k</b>	<b>L</b>			Mining						1449 (A. 1997)	2	<b>À</b>	<u> </u>	the		<u>k</u>	Marthan
and and a second states	atain her dia an	and the second	Section -	Andia .				. <u>k.</u>	L	1. 1	L		L										diara.	attanin	atter da.	in the second second		heristing.
						ll	NET TALK TO BE TO BE TO							and the second s	a star and a star	alte alte alte alte alte	4.0	La -dia ala ala ala	-1.0 -0.1 <u>an</u> 0.0 10		Labor Labor Johns		Looper Jacob	and any second s	and the sector		AND NOT DRAW	a store with

### Appendix B – Tools

The Colab file used for the energy regression part is:

https://colab.research.google.com/drive/1wbdo8p09oRSELITFzBZrL8zmzLuplh7J?usp=sharing

For the other tasks, other Colab files were used but one could simply replace the initial git clone link

It was mostly run with the L4 GPU, which gave the best efficiency compared to the cost. One can also run it on the free version, but importantly the paid version gives access to more CPU cores, that are used to increase the number of workers in the dataloader greatly improving the GPU utilization.

Free Colab Run T4 GPU – 2 CPU cores. CPU utilization is maxed out and GPU utilization is around 25%



Paid Colab Run L4 GPU – 12 CPU cores. CPU utilization is moderate and GPU utilization is around 75% for a beefier GPU



### Appendix B – Tools

In our Weights & Biases workspace we mainly tracked these following metrics:

Train\_loss: Training loss in the domain we are training

Val\_loss: Validation loss in the domain we are training

Val\_loss\_on\_ActualEnergy: Validation loss calculated after transforming prediction back to normal

energy val\_MSE\_on\_ActualEnergy val\_loss median\_train\_loss - FinGreaterZenith LR001 wZ - FinGreaterZenith LR001 wZ - FinGreaterZenith LR001 wZ - FinLessZenith LR001 wZ FinLessZenith LR001 wZ - FinLessZenith LR001 wZ OutputMLPHeadHiddenLayerLogEN OutputMLPHeadHiddenLayerLogEN OutputMLPHeadHiddenLayerLogEN OutputMLPHeadHiddenLayer OutputMLPHeadHiddenLayer OutputMLPHeadHiddenLayer •  $\mathbf{T}$  $\mathbf{T}$ 0.26 3.6e+10 0.28 0.24 3.4e+10 0.26 0.22 3.2e+10 0.24 0.2 3e+10 0.22 0.18 2.8e+10 0.2 0.16 2.6e+10 0.18 0.14 2.4e+10 0.16 trainer/global\_step ai**^**er/global\_step 0.12 100k 150k 200k 250k 300k 50k 50k 100k 150k 200k 250k 300k 50k 100k 150k 200k 250k 300k

### Appendix C – Neural network

The neural network was able to make simple but poor predictions. Not much time was spent on any form of optimization.

The goal of the NN was to overcome all obstacles related to data loading, running the code, and inference.

Gave great insight into the classes and data structures present in the code and data.

MLP1 - HP	Value
N_hiddenLayers	1
n_1	128
dropout	0.1

MLP2-HP	Value						
N_hiddenLayers	0						

### Appendix D – GNN

During data loading, all events are padded or truncated to an equal length.

The GNN model does not require equal-sized input data, so training with and without padding was performed.

Removing padding meant that the input data could not be collected in a single tensor forcing us to implement a for-loop over all events, seriously affecting the runtime of the algorithm.

Removing padding did, however, improve loss in the long run

### aggregated list = [] for b in range(batch size): batch x = x[b] # Shape: (seq len, feat dim) # Create mask for non-zero rows (valid events) valid mask = ~torch.all(batch x == 0, dim=1).to(device) valid x = batch x[valid mask] # Shape: (valid len, feat dim) edge index = self.get\_edge\_index(valid\_x).to(device) x\_ = self.convs(valid\_x, edge\_index) x\_ = torch.cat([x\_, valid\_x], dim=1).to(device) # Concatenate GNN output wit x\_ = self.mlp1(x\_) min pool = torch.min(x , dim=0)[0] max pool = torch.max(x , dim=0)[0] sum pool = torch.sum(x , dim=0) mean\_pool = torch.mean(x\_, dim=0) aggregated\_list.append(torch.cat([min\_pool, max\_pool, sum\_pool, mean\_pool]))

### Appendix D – GNN

Only 'small' runs with 1 training dataset were finalized, due to unexpected termination of runtimes for larger runs.

Without any error messages, resolving the issue proved difficult.

Due to time constraints and limited Colab tokens, the issues did not get resolved.

### GNN\_finalRun 🖉

Notes	What makes this run special? 🖉
Tags	+
Author	D david3435johansen
State	① Crashed
Start time	June 9th, 2025 9:49:39 AM
Duration	7h 49m 21s

### Appendix D – GNN

Due to long runtimes, only manual HP optimization has been performed. The parameters of the best run are shown in the table.

Hyper parameter	Value
K_neighbours	8
N_convs	4
Dropout	0.1
N_hidden_mlp1	1 (640)
N_hidden_mlp2	2 (64, 64)

### Appendix E – Transformer

We built our investigations of the transformer architecture for IceCube data regression on the work of Luc Voorhend. If one wants to learn more about or just use the model architecture, we encourage you to visit his github (<u>https://github.com/lucvoorend/IceCubeTransformer</u>) which gives a wonderful overview of the architecture. Additionally, our Githubs (Linked in the last slide of the presentation) show our changes to the model for the different use cases.

Here is an overview of the most important parts of the model we utilize also showing the structure of the files:

- ----- training\_and\_inference/ # Core scripts for the transformer model
- | |---- model.py # Script defining the model structure. (We used this to change the architecture)
- | |---- dataset.py # Script with Dataset Class (We used to change the domain and filtering process)
- ——— dataloader.py # Dataloader implementation (Used to make changes to the variables that were propagated to the dataset)
- $\vdash$ —— loss.py # Loss functions (Used to try out different loss functions)
- $\vdash$ —— train.py # Script to train the model

src/

- ├── inference.py # Script to run inference and evaluate the model
- └── config.yaml # Config file controlling settings for training and inference

### Appendix E – Transformer

class regression\_Transformer(nn.Module):

Regression transformer class:

- contains an input embedding layer, position embedding layer, transformer layers, and output layers
- applies the transformer to a sequence of embeddings
- returns a single predicted target value
- .....

self.input\_embedding = nn.Linear(input\_dim, embedding\_dim)
self.position\_embedding = nn.Embedding(seq\_dim, embedding\_dim)
self.layers = nn.ModuleList([EncoderBlock(embedding\_dim, n\_heads, dropout) for \_ in range(n\_layers)])
self.layer\_norm = nn.LayerNorm(embedding\_dim)
self.linear\_regression = nn.Linear(embedding\_dim, output\_dim)

def forward(self, x, target=None, event\_lengths=None, original\_event\_n\_doms=None):
 seq\_dim\_x = x.shape[1]
 device = x.device

input\_emb = self.input\_embedding(x).to(device)
pos\_emb = self.position\_embedding(torch.arange(seq\_dim\_x, device=device))
pos\_emb = pos\_emb.unsqueeze(0).expand(x.shape[0], -1, -1) # Shape: (batch\_size, seq\_dim, emb\_dim)

x = input\_emb + pos\_emb

# Feed the output of the transformer to the pooling layer batch\_dim, seq\_dim\_x, emb\_dim = x.shape[0], x.shape[1], x.shape[2]

# Aggregate x over the sequence dimension to the event length
row\_indices = torch.arange(seq\_dim\_x).view(1, -1, 1) # Shape: (1, seq\_dim, 1)
row\_indices = row\_indices.expand(batch\_dim, -1, emb\_dim)
row\_indices = row\_indices.to(device)

mask = row\_indices < event\_lengths.view(-1, 1, 1).to(device) # Shape: (batch\_size, seq\_dim, emb\_dim)</pre>

# Apply mask to x and do mean pooling x = x.masked\_fill(mask == 0, 0) x\_mean = x.sum(dim=1) / event\_lengths.view(-1, 1) # Shape: (batch\_size, emb\_dim) y\_pred = self.linear\_regression(x\_mean) On the left is the part of Luc's code defining the regression model. One can here rather easily follow what the model does which, for us, made the architecture very approachable and is why we include it explicitly here.

It starts with the embedding of both the simple input embedding and the position embedding being summed.

This is then input to the encoding layers which we usually repeat 4 times (n\_layers=4). Here, a few things are hidden under the hood namely the attention mechanism, FFN, as well as the LayerNorms and skip connections for better training and backpropagation.

After the encoding blocks the padded sequence elements are removed and mean aggregation is done.

Finally a linear layer does the final prediction mapping to the output dimension

### Appendix F – Transformer - Energy

Here we present our approach to doing energy regression using the Transformer architecture focusing on optimization steps and more fundamental changes we made to the model.

To start off we made a Github repository focusing on the energy regression. Here we made a model for energy regression by implementing and changing parts of Luc's code which is designed and hard coded for angular regression.

With the model running as intended we started the cruelling task of optimization. To optimize the transformer model for energy regression we just tried different configurations of hyperparameters (Grad student search).

With long training times, we did our best to do this in a structured manner focusing mostly on learning rate, sequence length, embedding dimension and batch size.

### Appendix F – Transformer - Energy Learning rate

With learning rate being as important as it is we tried a few different things with this. Namely, we tried using a dynamic learning rate, as shown on the right. The idea here is to let the model explore the loss landscape in the beginning with a later fine tuning.

However, we found that a simple constant learning rate would approach the general minimum validation loss we saw faster and more consistently. This was therefore used for the final models with values around 0.001.

It would be interesting to explore learning rate even more and maybe trying with a simple decaying learning rate.



### Appendix F – Transformer - Energy Loss function

We tried using a few different loss functions as seen on the right. Generally, we used MSE loss, but few runs used MAE. Huber loss was used for some initial trials in the true energy domain to try to mitigate the effects from the skewed values by evaluating errors on smaller values quadratically (MSE) and larger values linearly (MAE). We however quickly realized that the idea of using the true energy domain had to be abandoned.

```
def MSE_loss(y_pred, target):
    y_pred = y_pred.squeeze()
    target = target.squeeze()
    loss = target = resp((); prod = target
```

```
loss = torch.mean((y_pred - target) ** 2)
return loss
```

```
def MAE_loss(y_pred, target):
    y_pred = y_pred.squeeze()
    target = target.squeeze()
```

```
loss = torch.mean(torch.abs(y_pred - target))
return loss
```

```
def Huber_loss(y_pred, target, delta=0.1e6):
    y_pred = y_pred.squeeze()
    target = target.squeeze()
```

loss = F.huber\_loss(y\_pred, target, delta=delta, reduction='mean'
return loss

### Appendix F – Transformer - Energy Trials

Here we present a few of the later trials with their distinct hyperparameters. Generally, we see that they approach what we have denoted the 'Hand pattern'.







### Common HP

Parameter	Value
Batch Size	64
Num Layers	4
Num Heads	2
Dropout (FFN)	0.1

*There are naturally many more runs with other HP that we tried. These often also approach the 'Hand' or something worse.* 

Parameter	Value
Learning rate	Dynamic
Sequence length	256
Embedding Dim	256

Parameter	Value							
Learning rate	Dynamic							
Sequence length	512							
Embedding Dim	256							

Parameter	Value							
Learning rate	0.0005							
Sequence length	256							
Embedding Dim	128							

*\*Dynamic learning rate follows the curve on the earlier plot* 

### Appendix F – Transformer - Energy Investigating the outliers

As we have seen beyond multiple model architectures and HP parameter choices the energy predictions tend towards the mean. We turned our attention to these outliers to see whether the points are random or fall into some specific category.



As shown in the presentation, there is a clear correlation between predictive power and the zenith angle. We've seen that the model is not able to learn from particles going through the earth before hitting the detector. Additionally, we've seen a connection to the number of DOMs.



The plot on the left shows the average number of DOMs for the outliers and normal points for different energies. Interestingly, it is very constant for the outliers

For this reason, we tried adding the number of DOMs to the final linear layer thinking it could learn to understand this behavior better. This changed basically nothing, but there was no time to do optimization.

The next idea was to train the model on the particles coming from above and below the earth separately. The idea was to see whether the problem was the combination of the data. Also to see how good the model could learn without the clear outliers present.

### Appendix F – Transformer - Energy Zenith filter

To do the zenith specific training a zenith filter was introduced in the dataset class. For the middle cut at pi/2 presented in the slides, we saw that there is clearly little to no learning happening for the particles going through the earth with some learning for those that don't.



On the left is shown the plot for the model trained exclusively on particles above the horizon. Here, a colormap determined by the zenith angles shows that the horizontal line is dominated by higher zenith values. We therefore expect a more aggressive cut to get rid of this, but didn't have time to try this out

### Appendix F – Transformer - Energy Further Work



The scatter plot on the left is the result from an initial attempt at running the model in the log(E/N) domain with no optimization. One can maybe optimize this to work better than we've seen for the log(E) domain. More work can be focused on this.

The code below shows the NN we quickly tried at the last step before the output of the model. An initial test showed no improvement to earlier instances, but it would be interested to try working more with this.

A similar situation occurred for the Mean, Sum, Min and Max aggregation test.

```
self.output_mlp_head = nn.Sequential(
    nn.Linear(embedding_dim+1, (embedding_dim+1)*4),
    nn.ReLU(),
    nn.Dropout(dropout),
    nn.Linear((embedding_dim+1)*4, output_dim))
```

While grad student optimizing these runs, multiple different VMF loss functions were implemented. These included:

- Constant Penalty Terms
   Variable Penalty terms
   An openalty Terms
   Constant
- No Penalty Terms

Ultimately having no penalty terms yielded the best result as it seemed like the penalty needed tuning to not overrule the regular gradient descent. As this was very time consuming it was deemed outside the scope of the project.

\*This is likely also an artefact of not optimizing enough.

Also, we experimented with the learning rate of the models.

Many of the early models had a validation plot like the one on the right. This could potentially be due to a high learning rate; the models could not settle in a minimum but kept jumping out.

Unexpectedly the models consistently yielded better results for the checkpoint of the last epoch instead of the lowest loss epoch.



Predictions for the last epoch of this model were seen in the powerpoint. The lowest epoch can be seen below.



On the azimuth plot, the lowest blob is to be expected as the domain is periodic, and the predictions are shifted to match the physical interpretation. The middle blob, however, is because when the zenith angle has opposite sign, this corresponds to the azimuthal angle being rotated by  $\pi$ .

This effect tends to fade as the model continues running.

Work on attempting lower learning rates was initiated as the project was coming to its end. The final, yet intermediate, results can be seen below.

This run is split across two graphs due to technical difficulties. Both can be seen below. The left one contains the first 10 epochs. The right the last 25.



It is clear we do not see the same jumping behaviour as before, but the loss does not go lower (at least for our training time) than it did for a higher learning rate. It does however go steadily down which is why there is merit to this method. Potentially one could implement a learning rate that slowly decreases with epochs, yet it has to stay rather small from the beginning. The predictions for this are on the next slide.

Predictions for the low learning rate models.



The models have similar performance as the ones with higher learning rate; thus, no considerable improvement is made. The results did need more patching up, that is they originally emerged in different clusters due to the periodicity of the domain.

Playing with the learning rate is an interesting and obvious place for further work.

Note the interesting brick pattern on the left stemming from the transformations.

• We tried also to visualize the attention weights in the model training of angular regression. The relevant parameters for this model can be seen in the table.



All three plots stem from a model which has two attention heads and had 4 attention layers. On the left is a single head. In it we see clear vertical lines which are indications that certain specific keys (and therefore DOMs) are important. Some of the corresponding querys also light up a lot indicating these have strong connections to each other and that the model should take these into account when analyzing the keys. Also, the pattern only appears in a square in the top left corner. This is due to the specific event, for which the attention is visualized, only had a certain number of DOMs triggering, thus leaving the rest of the sequence length to be padded. Also, this padding is seen on the right-hand plot for which the actual event was even smaller.

On the right, we see both attention heads visualized at the same time. We see both have different attention patterns which agrees with the fact that different heads analyze (catch) different patterns. Here, the left one tends to connect events far (\*) from each other whereas the right-most one puts emphasis on the singular events themselves. Note also the squares on opposite sides (\*) light up.

(\*) We do not know which square in the attention weight plot corresponds to which DOM. Since the order is unknown to us, we cannot in principle relate this to the physical situation. Intuitively though, one could think that, say, pulses at opposite sides of the event would have increased value in determining the direction. One could then imagine a plot like the right-most one to represent this situation.