# Al Hands-On

By Andrii Tykhonov [andrii.tykhonov@cern.ch]

**Code:** https://gitlab.cern.ch/andrii/mlregressioncalo/-/tree/tutorial

ISAPP Lecce, Italy. June 12, 2025



# Preface: Installing Software















# Operating system: Linux or MacOS

Supported and tested operating systems:

**Linux** (currently tested on Ubuntu, but others like Debian, Fedora etc. should work just fine) lacksquare



Mac OS



Note on Windows — while it is not forbidden in the tutorial — it is neither tested nor fully supported, so it will be at your own risk. If you have a windows machine, it is adviced to install Linux either as a second operating system or in a virtual environment (e.g. through VirtualBox). Please contact me in advance if you have a Windows machine and never worked with Linux before.





# Software prerequisite: Miniconda is (almost) all you need! 4

### Install **Miniconda**:

- lacksquarelinux-installation
- lacksquare
- For example, in linux (basically same in MacOs but using curl instead of wget):

→ bash ~/Miniconda3-latest-Linux-x86 64.sh

- one paragraph to the setup script you can easily remove if you want to delete conda

→ source ~/.bashrc

(or ~/.zshrc - depending on which shell you are using)



### Follow the instructions in: https://www.anaconda.com/docs/getting-started/miniconda/install#macos-

Use **Terminal installer** (not graphical one): it allows to easily install/replace and experiment with Miniconda - everything will be placed in your home directory instead of the system one, so you will avoid potential conflicts with already pre-installed python versions etc.

→ wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86 64.sh

 $\rightarrow$  You will have to agree license agreement etc (press "q" to exit license agreement in the terminal  $\bigcirc$  ).

→ When it prompts "Choose an initialization options:" choose YES. If you are worried, you can make a backup of your profile initialization scripts (~/.bashrc or ~/.zshrc depending on the shell you use), but in principle all what conda does is adding











## Setup ML software: Tensorflow

Throughout the tutorial we will use both **Tensorflow** and **Pytorch** frameworks. We will install those in a two separate "environments" of conda and you will be able to easily switch between the two. You will appreciate the convenience and power of Conda – it allows you to install ML software (hopefully) quickly, (hopefully) gracefully, and without experiencing library conflicts etc. (unless something goes really wrong - but we are here to help you in this case ...). In the first part of the tutorial we will work with Tensorflow since (arguably) it is more simple/intuitive to use.

### Install **Tensorflow**

- 1. Add conda-forge channel to look for software:
- 2. Create new environment that we will call "tf":
- 3. Activate the "tf" environment:

**Comment on activation (step 3):** this has to be done in each new shell. If you want it to active the environment automatically at the start of a new shell, you can add the above command to your ~/.bashrc [or ~/.zshrc]

Note on the tensorflow versions: sometimes one needs to experiment with versions, for example I have experienced problems with latest tensorflow versions, so I downgraded to 2.16 on Ubuntu (2.15 on MacOS): conda create -n tf tensorflow=2.16









## Test Tensorflow #1 ...

- Activate the Tensorflow environment in conda (see previous page for mode details): conda activate tf
- See if Tensorflow libs are there and working: python
  - >>> import tensorflow as tf
  - >>> print (tf. version )
  - 2.16.XX

### NOTE USE Tensorflow version 2.15 or higher! ### If nothing crashes so far - things seem to work so far.. ### First import of tensorflow on some systems (e.g. MacOS) may take ### a while - be patient (it will cache and work faster afterwards)



### Test Tensorflow #2...

Let's run a mock-up model training:



• You should see something like:

Epoch	1/5								
32/32	[=======]	_	0s	206us/step	_	loss:	1.5475	_	accu
Epoch	2/5								
32/32	[========]	-	0s	158us/step	-	loss:	1.2575	_	accu
Epoch	3/5								
32/32	[======]	_	0s	136us/step	_	loss:	1.0151	_	accu
Epoch	4/5								
32/32	[======]	_	0s	129us/step	_	loss:	0.8074	_	accu
Epoch	5/5								
32/32	[=======]	-	0s	144us/step	-	loss:	0.6322	_	accu

>>> l = tf.keras.losses.SparseCategoricalCrossentropy(from\_logits=True) # we will learn about it later # train the model

racy: 0.0000e+00

racy: 0.0000e+00

uracy: 0.0000e+00

uracy: 0.0000e+00

uracy: 0.0000e+00





# If things go wrong ...



### **Possible solutions:**

Re-install another tenosorflow version (you may have to experiment with a few different versions):

```
conda remove -n tf --all
conda create -n tf tensorflow=2.16
```

Mask out your GPU\*:

```
export CUDA_VISIBLE_DEVICES=""
```

\* If you you have an Nvidia GPU and the the solution 1 (re-installing different tensorflow versions) do not help, try forcing tensorflow NOT to use the GPU (it is OK for the sake of this tutorial; in the future you may tweak your software setup to fully profit of your nice GPU hardware)

- these are few typical problems I encountered myself ... Unfortunately there might be more, but normally with the slight help of google, chatgpt, stackoverflow and a little prayer – things will work ;-)

```
# delete existing tf environment from conda
# install a specific tensorflow version (e.g. 2.16 ubuntu, or 2.15 for MacOs)
```







# Choose your weapon (text editor)

- It is perfectly fine to use your favorite python editor throughout the tutorial: vim, emacs, nano, eclipse+pydev, ...
- If you don't have one, for sake of simplicity and convenience it is suggested to install **jupyter**: conda activate tf conda install jupyter
- Create a directory for the code and run the editor there: mkdir mycode cd mycode jupyter notebook

# open a separate terminal window where you will run the editor







# Using jupyter as text editor

▼ - New	1 Upload
Modified	File Si

### • Jupyter will run a python text editor in your browser, there you can create new (click right mouse button)



# Using jupyter as text editor

- Let's call our file  ${\tt mycode.py}$  – double click on it and you will enter the editor

		localhost		S =	
☆ Start Page		C Home			🔵 Home
💭 Jupyter					
File View Settings Help					
Files O Running					
Open Download Rename Duplie	cate Move to Trash			▼ • New	± Upload C
■ /					
Name			-	Modified	File Size
V D mycode.py				now	0 В



# Using jupyter as text editor

• Et voilà!



you run the editor - jupyter, in the other one - the code itself):

### >>> python mycode.py >>> Tensorflow version 2.16.1

	ی چھ ا	
C Home		

• Keep a separate terminal window to run your code (essentially you have two terminal windows, in one of those



# Part I: Our first Neural Network



Adapted from: <u>https://www.tensorflow.org/tutorials/quickstart/beginner</u>

### 13

## Our first dataset: MNIST

### First we get some data that we want to train our NN ... Create a new file:

import tensorflow as tf
mnist = tf.keras.datasets.mnist
(x\_train, y\_train), (x\_test, y\_test) = mnist.load\_data()
x\_train, x\_test = x\_train / 255.0, x\_test / 255.0

### Let's have a look inside the data...

python -i first\_nn.py
>>> x\_train.shape
(60000, 28, 28) # - 60000 images of digits 28x28 pixel each
>>> y\_train.shape
(60000,) >>> # - 60000 labels corresponding to a number (from 0 to 9)

Reference: https://www.tensorflow.org/tutorials/quickstart/beginner





## Our first dataset: MNIST

### **MNIST dataset** (Modified National Institute of Standards and Technology database):

• A database of handwritten digits in the format of 28 x 28 pixel b/w images (pixel intensity) encoded in 1 byte, from 0 to 255)



(a) MNIST sample belonging to the digit '7'.

Reference: doi.org/10.3390/app9153169

(b) 100 samples from the MNIST training set.

Our goal – develop a NN that can classify a hand-written image telling wether it corresponds to 0, 1, ... or 9



## Our first dataset: MNIST

- conda install matplotlib
- Add visualization code to our first nn.py and run it:

```
import matplotlib.pyplot as plt
• • •
plt.imshow(data train[0], interpolation='none', cmap='gray')
plt.show()
plt.imshow(data train[1], interpolation='none', cmap='gray')
plt.show()
```



### • Let's get a habit of doing visualization/debugging of our data — we will need marplotlib library:

- # comment out this line afterwards # comment out this line afterwards # comment out this line afterwards # comment out this line afterwards





# On the Neural Networks (NN) in Tensorflow & Keras...

Tensorflow is normally used with **Keras wrapper API** (comes as a part of tensorflow installation) which allows to create and manipulate Neural Networks in intuitive way, composed of "layers" stacked one after another

**tf**.keras.layers.Dense(4,...)

tf.keras.layers.Input(3,..)





**tf**.keras.layers.Dense(1,...)



# Constructing our first model (NN) in Tensorflow

### Add NN model to our first nn.py:

```
model = tf.keras.models.Sequential([
  tf.keras.layers.Flatten(input_shape=(28, 28)), # First we convert input image into a flat array of numbers
  tf.keras.layers.Dense(128, activation='relu'), # Next, we add a layer of 128 neurons
  tf.keras.layers.Dropout(0.2),  # Dropout randomly removes 20% neurons from the above layer
tf.keras.layers.Dense(10)  # Final layer will correspond to 10 probabilities
])
```



Input image converted into a flat array

128 neurons (activation function Relu)

dropout - randomly removes 20% (0.2) of neurons from the above layer during training 10 outputs - will be converted into probabilities of an image representing 0,1,2,...,9



## Let's inspect our model (NN)

### >>> python -i first nn.py

>>> model.summary()

### Model: "sequential"

Layer (type)	Output Shape	Param #	
flatten ( <mark>Flatten</mark> )	(None, 784)	0	2
dense (Dense)	(None, 128)	100,480	
dropout ( <mark>Dropout</mark> )	(None, 128)	0	0
dense_1 (Dense)	(None, 10)	1,290/	?

**Total params:** 101,770 (397.54 KB) **Trainable params: 101,770** (397.54 KB) **Non-trainable params:** 0 (0.00 B)



### Remember neuron structure & parameters:





### Side note on activation functions ...



# Why do we need activation functions?

- ulletsimple linear combination of inputs!
- ulletprobabilities



And there are more (leaky ReLU, tanh etc.)

Primary reason for activation functions is to introduce a **non-linearity** in the model, otherwise the output will be a Activations are also required in classification tasks in order to convert continue outputs into contained [0,1]



21

# Example of NN without activation function



In this example model is not able to learn the data representation (identify two classes of orange and blue points) ...



# Example of NN without activation function





# Activation function (ReLU) added





### Now getting back to our code...

## What does our model do?

### Try yourself:

### python -i first nn.py >>> predictions = model(x\_train[:1]).numpy() # let's process a first training image with our model >>> print (predictions)

[[ 0.37128526 -0.10725151 0.18021962 0.10659367 0.20349179 0.00092683 -0.45577508 -0.23256837 0.02991931 -0.14916359]

# ... these do not look like probabilities # To interpret model output as probabilities we process the output with Softmax function:



### >>> predictions = tf.nn.softmax(predictions).numpy() >>> print (predictions)

[[0.11647741, 0.05115878, 0.09075072, 0.15573394, 0.10448843, 0.07079597, 0.05859897, 0.12428293, 0.14535967, 0.08235319]]

# ... now it looks like probabilities, but they are clearly wrong, because we need to train our model first!



## Adding loss function and compiling the model

### Add to first nn.py:

```
loss fn = tf.keras.losses.SparseCategoricalCrossentropy(from logits=True)
# example of loss calculation
predictions = model(x train[:1]).numpy()
print ('loss=',loss fn(y train[:1], predictions).numpy()) # comment me out later
model.compile(optimizer='adam',
              loss=loss fn,
              metrics=['accuracy'])
```

# Note that from logits=True informs the loss function that the model output is not reduced # to probabilities [0 to 1], hence the loss will apply softmax function to the model output

# comment me out later

# accuracy = N correct guesses / N total





## Note on the loss function

### Categorial cross-entropy:

$$\mathrm{Loss} = -\sum_{i=1}^{C} y_i \log(\hat{y}_i)$$

Cnumber of classes (10 in our case)

 $y_i$  – true label for class *i* (either 1 or 0)

 $\hat{y}_i$ – predicted label for class *i* (in the 0 to 1 range)

### Try yourself:

```
>>> from tensorflow import keras
>>> import numpy as np
>>> y true = np.array([1, 2])
>>> y_pred = np.array([[0.05, 0.95, 0], [0.1, 0.8, 0.1]])
>>> scce = keras.losses.SparseCategoricalCrossentropy()
>>> print (scce(y true, y pred))
```

The two examples correspond to identical cases, the only difference is in the format of the true labels!





# Training and testing the model

### Add training part to the first nn.py and run it:

history = model.fit(x\_train, y\_train, epochs=5)



### Add testing part:

model.evaluate(x\_test, y\_test, verbose=2)

.. you will get something like: 313/313 - 0s - loss: 0.0760 - accuracy: 0.9762

======] -	1s	339us/step -	loss:	0.3030	_	accuracy:	0.9123
=====] -	1s	333us/step -	loss:	0.1431	_	accuracy:	0.9577
	1s	464us/step -	loss:	0.1061	_	accuracy:	0.9679
=====] -	1s	327us/step -	loss:	0.0880	_	accuracy:	0.9728
==========] -	1s	338us/step -	loss:	0.0750	_	accuracv:	0.9768

### That's it - we trained our first NN model!



### Recap

```
import tensorflow as tf
#import matplotlib.pyplot as plt
# get the dataset
mnist = tf.keras.datasets.mnist
(x train, y train), (x test, y test) = mnist.load data()
x train, x test = x train / 255.0, x test / 255.0
# example of input data
#plt.imshow(x train[0], interpolation='none', cmap='gray') # comment me out later
#plt.show()
#plt.imshow(x train[1], interpolation='none', cmap='gray') # comment me out later
#plt.show()
# create the model
model = tf.keras.models.Sequential([
  tf.keras.layers.Flatten(input shape=(28, 28)), # First we convert input image into a flat array of numbers
  tf.keras.layers.Dense(128, activation='relu'), # Next, we add a layer of 128 neurons
  tf.keras.layers.Dropout(0.2),  # Dropout randomly removes 20% neurons from the above layer
  tf.keras.layers.Dense(10)
])
# example of how model process the intput data
#predictions = model(x train[:1]).numpy()
#print (predictions)
# define the loss function
loss fn = tf.keras.losses.SparseCategoricalCrossentropy(from logits=True)
# example of loss calculation
#predictions = model(x train[:1]).numpy()
#print ('loss=',loss fn(y train[:1], predictions).numpy())
# compile the model
model.compile(optimizer='adam',
              loss=loss fn,
              metrics=['accuracy']) # accuracy = N correct guesses / N total
# train the model
history = model.fit(x train, y train, epochs=5)
# test the model
model.evaluate(x_test, y_test, verbose=2)
```

# comment me out later # comment me out later

# Final layer will correspond to 10 probabilities

#	comment	me	out	later
#	comment	me	out	later
#	comment	me	out	later

### first nn.py



## Recap (remove non-essential commented parts)

```
import tensorflow as tf
# get the dataset
mnist = tf.keras.datasets.mnist
(x_train, y_train), (x_test, y_test) = mnist.load_data()
x train, x test = x train / 255.0, x test / 255.0
# create the model
model = tf.keras.models.Sequential([
  tf.keras.layers.Dense(128, activation='relu'), # Next, we add a layer of 128 neurons
  tf.keras.layers.Dense(10)
])
# define the loss function
loss fn = tf.keras.losses.SparseCategoricalCrossentropy(from logits=True)
# compile the model
model.compile(optimizer='adam',
              loss=loss fn,
              metrics=['accuracy'])
# train the model
history = model.fit(x train, y train, epochs=5)
# test the model
model.evaluate(x_test, y_test, verbose=2)
```

tf.keras.layers.Flatten(input shape=(28, 28)), # First we convert input image into a flat array of numbers tf.keras.layers.Dropout(0.2), # Dropout randomly removes 20% neurons from the above layer # Final layer will correspond to 10 probabilities



## Recap (remove non-essential commented parts)

```
import tensorflow as tf
```

```
# get the dataset
mnist = tf.keras.datasets.mnist
(x_train, y_train), (x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0
```

```
# create the model
model = tf.keras.models.Sequential([
  tf.keras.layers.Dense(128, activation='relu'), # Next, we add a layer of 128 neurons
  tf.keras.layers.Dense(10)
```

```
# define the loss function
loss_fn = tf.keras.losses.SparseCategoricalCrossentropy(from logits=True)
# compile the model
model.compile(optimizer='adam',
              loss=loss fn,
              metrics=['accuracy'])
# train the model
history = model.fit(x train, y train, epochs=5)
# test the model
model.evaluate(x_test, y_test, verbose=2)
```

This is the core part of tensorlfow (or any other ML framework) – when the model is built, tensorflow becomes aware of its trainable parameters (wights) and bias values of every neuron in every layer)

tf.keras.layers.Flatten(input shape=(28, 28)), # First we convert input image into a flat array of numbers tf.keras.layers.Dropout(0.2), # Dropout randomly removes 20% neurons from the above layer # Final layer will correspond to 10 probabilities



## Recap (remove non-essential commented parts)

```
import tensorflow as tf
```

```
# get the dataset
mnist = tf.keras.datasets.mnist
(x_train, y_train), (x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0
```

```
# create the model
model = tf.keras.models.Sequential([
  tf.keras.layers.Dense(128, activation='relu'), # Next, we add a layer of 128 neurons
  tf.keras.layers.Dropout(0.2),
  tf.keras.layers.Dense(10)
```

```
# define the loss function
loss fn = tf.keras.losses.SparseCategoricalCrossentropy(from logits=True)
# compile the model
model.compile(optimizer='adam',
              loss=loss fn,
              metrics=['accuracy'])
```

```
# train the model
history = model.fit(x train, y train, epochs=5)
```

```
# test the model
model.evaluate(x_test, y_test, verbose=2)
```

This is the core part of tensorlfow (or any other ML framework) – when the model is built, tensorflow becomes aware of its trainable parameters (wights) and bias values of every neuron in every layers)

tf.keras.layers.Flatten(input\_shape=(28, 28)), # First we convert input image into a flat array of numbers # Dropout randomly removes 20% neurons from the above layer # Final layer will correspond to 10 probabilities

The partial derivatives w.r.t. trainable parameters are computed during the execution of "fit" given the input data. The parameters are updated in every "epoch" following the gradient descent method (e.g. "adam" is one of the most common types of gradient descent algorithms)





# Fun part: let's run some predictions!

First, remember that our model produces a set of 10 numbers as the input, however they are not constrained to [0,1] range since the constraining part (Softmax activation) was included in the loss function, but not in the model itself.

Hence, we need to convert our 10 number into 10 probabilities, in first nn.py add:

```
probability model = tf.keras.Sequential([
 model,
 tf.keras.layers.Softmax()
```

Now let's run the predictions:

```
python -i first nn.py
>>> import numpy as np
>>> import matplotlib.pyplot as plt
>>> plt.imshow(x test[17], interpolation='none')
>>> plt.show()
>>> model guss = probability model(x test[17:18])
>>> print ("The guessed number is:", np.argmax(model guss))
```

Try with the other images and see for yourself how good (or maybe not) our model is? :-) Don't close the console with the first nn.py yet! (See next slide...)











# Visualizing the training process

### **Don't close your console with** first nn.py yet, let's examine the history object:

- >>> print (history.history)
- {'loss': [0.29248857498168945, 0.1410820633172989, 0.10512221604585648, 0.08767301589250565,
- 0.9729166626930237, 0.9773499965667725]

### Now let's do some plots:

- >>> plt.plot(history.history['loss'])
- >>> plt.plot(history.history['accuracy'])
- >>> plt.legend(['loss', 'accuracy'])
- >>> plt.show()

0.07357344776391983], 'accuracy': [0.914900004863739, 0.9583333134651184, 0.9677833318710327,





## Let's complicate things a bit ...

In your first nn.py modify the fit function to the following and re-run the code:

history = model.fit(x train, y train, epochs=50, validation split=0.1)

**Now let's inspect the** history **object again and plot the numbers there**:

>>> print (history.history.keys()) dict keys(['loss', 'accuracy', 'val loss', 'val accuracy'])

- >>> plt.plot(history.history['loss'])
- >>> plt.plot(history.history['val loss'])
- >>> plt.plot(history.history['accuracy'])
- >>> plt.plot(history.history['val accuracy'])
- >>> plt.legend(['loss','val\_loss','accuracy','val\_accuracy'])
- >>> plt.show()




## Let's complicate things a bit (adding validation)...

In your first\_nn.py modify the fit function to the following and re-run the code:

history = model.fit(x\_train, y\_train, epochs=50, validation\_split=0.1)

Now let's inspect the history object again and plot the numbers there:

>>> print (history.history.keys())
dict\_keys(['loss', 'accuracy', 'val\_loss', 'val\_accuracy'])

- >>> plt.plot(history.history['loss'])
- >>> plt.plot(history.history['val\_loss'])
- >>> plt.plot(history.history['accuracy'])
- >>> plt.plot(history.history['val\_accuracy'])
- >>> plt.legend(['loss','val\_loss','accuracy','val\_accuracy'])
- >>> plt.show()

Try to experiment a bit, consider modifying model "horizontally" or "vertically" (adding/removing layers, number of neurons etc.). Can we arrive with a better (more accurate model) that, at the same time, does not overfit? ...





### Put our code in order

For convenience of further work, let's re-structure our code a bit:

cp first nn.py run fit.py

In run fit.py:

from model nn import model

• • • #model = tf.keras.models.Sequential([...]) # <-- remove or comment out this part</pre>

**Create an empty** init .py **along with the following** model nn.py:

```
import tensorflow as tf
model = tf.keras.models.Sequential([
  tf.keras.layers.Flatten(input shape=(28, 28)),
  tf.keras.layers.Dense(128, activation='relu'),
  tf.keras.layers.Dropout(0.2),
  tf.keras.layers.Dense(10)
])
```

### # <-- add this in the beginning

We didn't do anything new, we just restructured the code for easier manipulation with difference models. If you run run fit.py you will get the same result as with first\_nn.py



## Another way of defining Tensorlfow (TF) model

Let's rewrite our NN model in slightly different format (it will get useful for us later ...). Create model nn new.py:

```
import tensorflow as tf
input = tf.keras.Input(shape=(28, 28))
layer = tf.keras.layers.Flatten()(input)
layer = tf.keras.layers.Dense(128, activation='relu')(layer)
layer = tf.keras.layers.Dropout(0.2)(layer)
layer = tf.keras.layers.Dense(10) (layer)
model = tf.keras.Model(inputs=input,outputs=layer)
```

• In the run fit.py file: replace model nn with model nn new and run the run fit.py again ...

While nothing really changed in what the code does, the new format/definition of the model will come in hand while working with CNNs, in particular for understanding the dimensionality of tensors etc.



## Part II: Convolutional Neural Networks



### What are Convolutional Neural Networks?



Convolution is essentially a filter that slides through the image producing one output value per each position

In the example of 5x5 convolution filter, at every position the output value will be:

$$O(i,j) = \sum_{m=0}^4 \sum_{n=0}^4 K(m,n) \cdot I(i+m,j+n) \ + b$$

O(i,j)- output value at filter position (i,j) - in this example from 0 to 23 K(m,n)- the convolutional kernel value at row *m* and column *n* (from 0 to 4) I(i+m, j+n) – the corresponding pixel value being convolved – bias term  $\boldsymbol{b}$ 





### What are Convolutional Neural Networks?

With the 5x5 filter, the output will be another picture of dimension 24 x 24. Now imagine you have two such filters, so that the two output images will be produced:







### What are Convolutional Neural Networks?

With the 5x5 filter, the output will be another picture of dimension 24 x 24. Now imagine you have two such filters, so that the two output images will be produced:



- We apply convolutional layers one after another allowing the model to learn some deep features in the data ullet
- The 2D dimensionality necessarily decreases layer after layer
- Normally after the last layer will will end up with 1x1xN images, i.e. N numbers which we can then process with the usual NN

• The 3<sup>rd</sup> dimension equals the number of Conv filters at every step (usually is set to increase layer after layer, but not necessarily)



### CNN with our MNIST example

### Let's get back our number classification example: copy model nn new.py to model cnn.py and let's do some changes there:

```
import tensorflow as tf
```

```
#input = tf.keras.Input(shape=(28,28))
```

#layer = tf.keras.layers.Flatten()(input) # < -- we don't need to flatten at this point...

# convolutional part layer = tf.keras.layers.Conv2D(32,(4,4),strides=(4,4), activation='relu')(input) layer = tf.keras.layers.Conv2D(64, (7,7), strides=(1,1), activation='relu')(layer) layer = tf.keras.layers.Flatten() (layer)

```
# usual NN part remains the same
layer = tf.keras.layers.Dense(128, activation='relu')(layer)
layer = tf.keras.layers.Dropout(0.2)(layer)
layer = tf.keras.layers.Dense(10) (layer)
```

model = tf.keras.Model(inputs=input,outputs=layer)



input = tf.keras.Input(shape=(28,28,1)) # <-- third dimension specify number of channels in the image # <-- for example, b/w has only 1 channel, colored has 3 (RGB)



### CNN with our MNIST example

### Let's get back our number classification example: copy model nn new.py rewrite model cnn.py and let's do some changes:

```
import tensorflow as tf
```

```
#input = tf.keras.Input(shape=(28,28))
#layer = tf.keras.layers.Flatten()(input_) # <-- we don't need to flatten at this point...</pre>
# convolutional part
layer = tf.keras.layers.Conv2D(32,(4,4),strides=(4,4), activation='relu')(input)
print('Conv layer 1 shape:', layer.shape)
layer = tf.keras.layers.Conv2D(64, (7,7), strides=(1,1), activation='relu')(layer)
print('Conv layer 2 shape:', layer.shape)
layer = tf.keras.layers.Flatten() (layer)
# usual NN part remains the same
layer = tf.keras.layers.Dense(128, activation='relu')(layer)
layer = tf.keras.layers.Dropout(0.2)(layer)
```

```
layer = tf.keras.layers.Dense(10) (layer)
```

model = tf.keras.Model(inputs=input,outputs=layer)

python model cnn.py

Check the output for yourself and let's see if we understand it ...





### CNN with our MNIST example

### **Open** run\_fit.py **and replace the NN with the CNN; re-run the training**:

#from model\_nn\_new import model # comment out the NN
from model\_cnn import model # import CNN instead

python run\_fit.py

Epoch 1/50						
1688/1688 [=================]]	_	1s	789us/step	—	loss:	0.28
Epoch 2/50						
1688/1688 [=================]]	_	1s	765us/step	—	loss:	0.11
Epoch 3/50						
1688/1688 [=================]]	_	1s	761us/step	—	loss:	0.07
Epoch 4/50						
1688/1688 [=============]]	—	1s	762us/step	-	loss:	0.06

Comped to NN model, the CNN model appears to have similar performance (check the loss and accuracy metrics of both models)...

However, keep in mind that we deal with very small images (28 x 28) — things will be different if we consider e.g. Megapixel scale pictures...

```
831 - accuracy: 0.9139 - val_loss: 0.1069 - val_accuracy: 0.9710
134 - accuracy: 0.9660 - val_loss: 0.0829 - val_accuracy: 0.9768
798 - accuracy: 0.9753 - val_loss: 0.0782 - val_accuracy: 0.9765
624 - accuracy: 0.9807 - val_loss: 0.0663 - val_accuracy: 0.9817
```





### Inspect our first CNN model

python -i model\_cnn.py
>>> model.summary()

Model: "functional"

Layer (type)	Output Shape	Param #	
<pre>input_layer (InputLayer)</pre>	(None, 28, 28, 1)	0	
conv2d (Conv2D)	(None, 7, <u>7</u> , 32)	544	?
conv2d_1 (Conv2D)	(None, 1, 1, 64)	100,416	?
flatten ( <mark>Flatten</mark> )	(None, 64)	0	
dense (Dense)	(None, 128)	8,320	?
dropout ( <mark>Dropout</mark> )	(None, 128)	0	
dense_1 (Dense)	(None, 10)	1,290	?

Total params: 110,570 (431.91 KB) Trainable params: 110,570 (431.91 KB) Non-trainable params: 0 (0.00 B)

Do the math yourself, do you arrive to the same numbers?





### On the power of CNNs...

Consider 4K image classification(~10 megapixels):





• Take a simple NN (not CNN) model architecture from the MNIST example (with  $\sim 100$  neurons in the first layer) – how many trainable parameters do you get? Compare it with the number of trainable parameters we had in our NN and CNN models (~100'000). The number of parameters in NN model looks big, right? (Comparable to GPT-2!)

• CNNs, on the other hand, allow to extract features from the images of arbitrary size and resolution. The complexity of CNN model is defined only by the filter dimension, their number, and amount of layers – all these can be set to whatever number depending on how complex is the model that you want to build (how much trainable date you have etc.)



# CNN example with the DAMPE cosmic ray detector



### DAMPE example

### Predicting cosmic ray (or gamma ray) particle direction from a signal (shower) in DAMPE calorimeter:





"Heart" of the DAMPE detector - the BGO imaging calorimeter





### DAMPE example

### Predicting cosmic ray (or gamma ray) particle direction from a signal (shower) in DAMPE calorimeter: DAMPE detects particle spatial information and energy deposition in two orthogonal projections (XZ and YZ)



DAMPE XZ E=1.416 TeV

DAMPE YZ E=1.416 TeV



## DAMPE example: classical approach (no AI)

Predicting cosmic ray (or gamma ray) particle direction from a signal (shower) in DAMPE calorimeter: DAMPE detects particle spatial information and energy deposition in two orthogonal projections (XZ and YZ)



In a classical approach, the particle trajectory in the calorimeter is estimated by performing a linear regression fit of the line through the points in calorimeter, where the contribution of every point in the fit is weighted by the energy deposition in this point



### DAMPE example: CNN

- We will use CNN instead of the classical linear regression to predict particle direction
- We will train it on a sample of ~140'000 simulated particle showers in DAMPE



n - to predict particle direction le showers in DAMPE



## DAMPE example: getting the data

### • Download and setup the training DAMPE data and the corresponding software:

git clone https://gitlab.cern.ch/andrii/mlregressioncalo.git cd mlregressioncalo git switch tutorial # do it in every new console (cd [path-/mlregressioncalo]; source ./setup.sh) source ./setup.sh

• Inspect the content of DAMPE package & data



*Do yourself - check the dimensions of the data arrays:* data['caloimages'].shape *etc.* 

The same 4 numbers but obtained with the standard linear regression algorithm — we will not use it for training, but will keep it for reference when comparing with the CNN predictions





### DAMPE example: visually examining the data

• Let's inspect some of the data events:

>>> plot\_dampe\_event(data['caloimages'][10],data['truthdata'][10])





### DAMPE example: visually examining the data

Let's inspect some of the data events: 



You can experiment with the different events in the dataset, for example in the event 19 one can see a very distinct difference between the truth and reconstructed particle direction.

### >>> plot dampe event(data['caloimages'][10],data['truthdata'][10],data['standardrecdata'][10])





### DAMPE CNN training ...

### • Let's get back to our CNN training code from the MNIST example, we will use it as a base:

```
cp run fit.py run fit dampe.py
cp model_cnn.py model_cnn_dampe.py
```

### • In run fit dampe.py change:

```
# before
#from model cnn import model
from model cnn dampe import model
                                     # now
from dampe import get_dampe_data
from pickle import dump
• • •
# get the dataset
#mnist = tf.keras.datasets.mnist
#(x_train, y_train), (x_test, y_test) = mnist.load_data()
#x train, x test = x train / 255.0, x test / 255.0
data = get dampe data()
x, y = data['caloimages'], data['truthdata']
                                             #
• • •
```

```
# before
# now
#
```

```
#loss fn = tf.keras.losses.SparseCategoricalCrossentropy(from logits=True)
loss fn = tf.keras.losses.MeanAbsoluteError()
model.compile(optimizer='adam',loss=loss_fn, metrics=['mean_squared_error']) #metrics=[`accuracy'])
history = model.fit(x, y, (epochs=5) validation split=0.1) # training part remains the same
# After the fitting part let's save the model and history for further analysis
model.save_weights('model_cnn_dampe.weights.h5')
dump(history.history, open('history cnn dampe.p','wb'))
# comment out the rest of the code afterwards (plt.plot(history.history ... etc.)
```

```
# before
# now
```

Note that we reduce number of epochs since training DAMPE models is more time consuming ...



### DAMPE CNN model ...

### • Now we need to modify the model itself, keeping the changes minimal. In the model cnn dampe.py do the following:

```
#input = tf.keras.Input(shape=(28,28,1)) # before
input = tf.keras.Input(shape=(14,22,1)) # now
layer = tf.keras.layers.ZeroPadding2D((((1,1),(1,1)))(input )
```

```
# convolutional part
#layer = tf.keras.layers.Conv2D(32,(4,4),strides=(4,4), activation="relu")(input)# before
layer = tf.keras.layers.Conv2D(32, (4, 4), strides=(4, 4), activation="relu")(layer) # now
#layer = tf.keras.layers.Conv2D(64, (7,7), strides=(1,1), activation="relu")(layer) # before
layer = tf.keras.layers.Conv2D(64,(6,4),strides=(1,1), activation="relu")(layer)
layer = tf.keras.layers.Flatten()(layer)
```

```
# NN part
layer = tf.keras.layers.Dense(128, activation="relu")(layer)
#layer = tf.keras.layers.Dropout(0.2)(layer)
#layer = tf.keras.layers.Dense(10)(layer)
layer = tf.keras.layers.Dense(4) (layer)
```

```
# this remains the same as before ...
model = tf.keras.Model(inputs=input,outputs=layer)
```

• We are ready to run the DAMPE model training:

python run fit dampe.py

```
• • •
Epoch 1/50
3993/3993 [=============================] - 3s 684us/step - loss: 76.2474 -
34.6831 - val_mean_squared error: 2725.1121
• • •
```

Test/print the layer shapes/dimensions at different steps, let's see if we understand it ... # now # same as before (not changed)

# same as before (not changed) # comment out the dropout part # before # now (instead of 10 neurons, now we have 4)

mean squared error: 11470.5166 - val loss:

Note the loss value at the first iteration (we will compare it later with the loss of the other model)









## DAMPE CNN model from the paper

### • Let's try a <u>deeper model (from the paper)</u>, cp model cnn.py model cnn dampe paper.py and modify it:

```
#input = tf.keras.Input(shape=(28,28,1)) # before
input = tf.keras.Input(shape=(14,22,1)) # now
```

```
# convolutional part
#layer = tf.keras.layers.Conv2D(32,(4,4),strides=(4,4), activation='relu')(input) # before
#layer = tf.keras.layers.Conv2D(64,(7,7),strides=(1,1), activation='relu')(layer)
layer = tf.keras.layers.Conv2D(128, (4, 4), activation="relu")(input)
                                                                                     # now
layer = tf.keras.layers.Conv2D(64, (4, 4), activation="relu")(layer)
layer = tf.keras.layers.Conv2D(32, (4, 4), activation="relu") (layer)
# ... what is the ouput shape after this layer? (see for yourself with print...)
layer = tf.keras.layers.Conv2D(100, (5,13), activation="relu")(layer)
# ... why is there (5,13) filter size used at this point? What is the output shape?
layer = tf.keras.layers.Flatten() (layer)
```

```
# NN part
#layer = tf.keras.layers.Dense(128, activation='relu'
#layer = tf.keras.layers.Dropout(0.2)(layer)
#layer = tf.keras.layers.Dense(10)(layer)
layer = tf.keras.layers.Dense(50, activation="relu")(
layer = tf.keras.layers.Dense(4, activation="linear")
```

```
# this remains the same as before ...
model = tf.keras.Model(inputs=input,outputs=layer)
```

)(layer)	# before		
lavor)	# # #		Model from:
(layer)	# 110 W #	Astroparticle Physics 146 (2023) 102795	
		Contents lists available at ScienceDirect Astroparticle Physics	ASTRO

**ELSEVIER** 

A deep learning method for the trajectory reconstruction of cosmic rays with the DAMPE mission

journal homepage: www.elsevier.com/locate/astropartpl

Andrii Tykhonov<sup>a,\*</sup>, Andrii Kotenko<sup>a</sup>, Paul Coppin<sup>a</sup>, Maksym Deliyergiyev<sup>a</sup>, David Droz<sup>a</sup>, Jennifer Maria Frieden<sup>b</sup>, Chiara Perrina<sup>b</sup>, Enzo Putti-Garcia<sup>a</sup>, Arshia Ruina<sup>a</sup>, Mikhail Stolpovskiy<sup>a</sup>, Xin Wu<sup>a</sup>

<sup>a</sup> Department of Nuclear and Particle Physics, University of Geneva, CH-1211, Switzerland

<sup>b</sup> Institute of Physics, Ecole Polytechnique Fédérale de Lausanne (EPFL), CH-1015, Lausanne, Switzerland







## DAMPE CNN model from the paper

### • In run fit dampe.py modify the CNN model import and run the script again:

#from model cnn dampe import model # before from model cnn dampe paper import model # now

# save model under different name #model.save weights('model cnn dampe.weights.h5') #dump(history.history, open('history cnn dampe.p','wb')) model.save weights('model cnn dampe paper.weights.h5') dump(history.history, open('history cnn dampe paper.p','wb'))

python run\_fit\_dampe.py

```
• • •
Epoch 1/50
3993/3993 [===========================] - 44s 11ms/step - loss:
val mean squared error: 229.4982
• • •
```

It will take 1-2 hours on conventional hardware to train this model... We can't wait that long, hence you can find this saved model and training history inside the dampe package downloaded earlier, see the saved models folder



Note that the loss value after the first iteration is considerably lower that with the first (simple) CNN model that we tried... Also, training takes considerably longer due to a more complex (dense) CNN model!





### Comparing the two CNN models

python

```
from pickle import load
import matplotlib.pyplot as plt
history1 = load(open('history_cnn_dampe.p', 'rb'))
history2 = load(open('history cnn dampe paper.p', 'rb'))
plt.plot(history1['loss'])
plt.plot(history1['val loss'])
plt.plot(history2['loss'])
plt.plot(history2['val loss'])
plt.show()
```





plt.legend(['loss (model simple)', 'val\_loss (model simple)', 'loss (model paper)', 'val\_loss (model paper)'])





## Fun part: inference/prediction with the DAMPE CNN

# This is test dampe images.py file for DAMPE CNN testing

```
from dampe import get_dampe_data, plot_dampe_event
import tensorflow as tf
from model cnn dampe import model
```

data = get dampe data()

```
# plot the results with the standard algorithm (linear regression)
standardprediction = data['standardrecdata'][10]
plot_dampe_event(data['caloimages'][10], data['truthdata'][10], standardprediction)
```

```
# plot the results with CNN
model.compile()
model.load weights('model cnn dampe paper.weights.h5')
prediction = model(data['caloimages'][10:11])
prediction = prediction[0] # prediction is done in batches. We have a "batch" of 1 event
plot_dampe_event(data['caloimages'][10], data['truthdata'][10], prediction)
```

```
# caclcualte yourself accucay of standard approach VS CNN
import numpy as np
```

```
n = 1000
```

print ("Standard algorithm mean absolute error:",

```
np.sum(np.abs(data['truthdata'][:n]-data['standardrecdata'][:n])) / (4*n))
prediction = model(data['caloimages'][:n])
print ("CNN mean absolute error:",
```

np.sum(np.abs(data['truthdata'][:n]-prediction)) / (4\*n))

• Let's compare performance of classical DAMPE algorithm and the CNN one. Create test dampe images.py and run it:











## Recap of the DAMPE CNN model (training script)

run fit dampe.py:

```
import tensorflow as tf
import matplotlib.pyplot as plt
from model cnn dampe paper import model
from dampe import get dampe data
from pickle import dump
```

```
data = get dampe data()
x, y = data['caloimages'], data['truthdata']
```

```
loss fn = tf.keras.losses.MeanAbsoluteError()
model.compile(optimizer='adam', loss=loss fn,
metrics=['mean squared error'])
history = model.fit(x, y, epochs=5, validation split=0.1)
```

```
model.save weights('./model cnn dampe paper.weights.h5')
dump(history.history, open('history cnn dampe paper.p','wb'))
```



## Recap of the DAMPE CNN model (model definition)

model cnn dampe\_paper.py:

import tensorflow as tf

input = tf.keras.Input(shape=(14,22,1))

# convolutional part layer = tf.keras.layers.Conv2D(128, (4, 4), activation="relu")(input ) layer = tf.keras.layers.Conv2D(64, (4, 4), activation="relu")(layer) layer = tf.keras.layers.Conv2D(32, (4, 4), activation="relu")(layer) layer = tf.keras.layers.Conv2D(100, (5,13), activation="relu")(layer) layer = tf.keras.layers.Flatten() (layer)

```
# fully-connected part
layer = tf.keras.layers.Dense(50, activation="relu")(layer)
layer = tf.keras.layers.Dense(4, activation="linear")(layer)
```

```
model = tf.keras.Model(inputs=input,outputs=layer)
```



## Exercise (20 mins)

- **Develop your own NN (not CNN) model for DAMPE** (similar to MNIST example discussed earlier)
  - $\rightarrow$  Experiment with architecture (number of layers, neurons, etc...)
  - $\rightarrow$  The goal is to try to obtain a simple NN model comparable in performance with the CNN
    - For simplicity, as a metrics of "performance" just use the value of loss (the lower the better)

first few iterations how good the model is...)

Modify the original CNN according to your intuition and try to beat the original DAMPE CNN!

No need to run the entire model training and visualization, just run a training with 2-3 iterations at most and see if your NN model converges as fast (with the similar or lower loss values) as the CNN model (usually it is already seen from the



Before we move on: let's repeat the DAMPE exercise while learning a bit of **PyTorch** (We will need PyTorch for the following part of the course)





### Install PyTorch

To avoid possible conflicts between the two frameworks, we will install PyTorch in a separate conda environment

- Install **PyTorch** 
  - 1. Create new environment that we will call "tr":
  - 2. Activate the "tf" environment:
- Test **PyTorch**

python

- >>> import torch as tr
- >>> print (tr.\_\_version\_\_)





## Test PyTorch NN training

• Run the following code snippet in python (let's call it test pytorch minimal.py):

```
import torch as tr
import numpy as np
model = tr.nn.Sequential( tr.nn.Linear(10, 10, dtype=tr.float64))
criterion = tr.nn.CrossEntropyLoss()
optimizer = tr.optim.Adam(model.parameters())
# one training step
y \text{ pred} = \text{model}(x)
loss = criterion(y pred,y)
loss.backward()
optimizer.step()
print ("Loss:",loss.item())
```

• If everything is installed correctly, you should see the log of the training:

Loss: 1.930114470175824

x = tr.tensor(np.random.rand(1000,10),dtype=tr.float64) # random sample of 1000 sets of numbers y = tr.tensor(np.random.rand(1000,),dtype=tr.int64) # random sample of 1000 set 'labels'





## PyTorch DAMPE example - training script

Let's rewrite our DAMPE CNN code in PyTorch:

cp run fit dampe.py run fit dampe torch.py, edit run fit dampe torch.py:

```
#import tensorflow as tf
import torch as tr
#from model cnn dampe paper import model
from model cnn dampe paper torch import model
```

#... after you obtained the dampe data in the usual way, convert it to torch tensors x, y = tr.tensor(x,dtype=tr.float32), tr.tensor(y,dtype=tr.float32)

```
#... we need to convert data format from [N,H,W,C] (Tensorflow) to [N,W,H,C] (Pytorch)
x = x.squeeze(3).unsqueeze(1)
#... you can check later by inserting print(x.shape) statement before and after...
```

```
#loss fn = tf.keras.losses.MeanAbsoluteError()
#model.compile(optimizer='adam', loss=loss fn, metrics=['mean squared error'])
criterion = tr.nn.LlLoss()
optimizer = tr.optim.Adam(model.parameters())
```

# continued on the next page ...

# mean absolute error loss



## PyTorch DAMPE example - training script

### ... edit run fit dampe torch.py continued:

```
#history = model.fit(x, y, epochs=5, validation split=0.1)
model.train() # set model in the training mode
history = { 'loss':[], 'loss_val':[] }
n total, n split, n batch = x.shape[0], int(x.shape[0] * 0.9), 32
for i in range(5): # loop over training epochs
    loss_train, loss_val, n_batches_train, n_batches_val = 0, 0, 0, 0
    for j in range(0, n_split, n_batch): # loop over batches in training sub-sample
        print (f"Processed: {j*100./n split:3.1f}%",end="\r") # status bar
        optimizer.zero grad()
        y pred = model(x[j:j+n batch])
        loss = criterion(y pred,y[j:j+n batch])
        loss.backward()
        optimizer.step()
        loss train+=loss.item()
        n batches train+=1
    model.valid() # set model in validation mode
        y_pred = model(x[j:j+n_batch])
        loss = criterion(y_pred,y[j:j+n_batch])
        loss val+=loss.item()
        n batches val+=1
    loss_train, loss_val = loss_train / n_batches_train, loss_val / n_batches_val
    print ("loss:",loss train, "loss val:",loss val)
#model.save weights('./model cnn dampe paper.weights.h5')
#dump(history.history, open('history cnn dampe smallmodel.p','wb'))
```





### we have that...

### **PyTorch DAMPE example - model definition**

### Let's rewrite the **model** in PyTorch,

```
cp model cnn dampe paper.py model cnn dampe paper torch.py,
edit model cnn dampe paper torch.py:
```

```
#import tensorflow as tf
import torch as tr
import numpy as np
```

#input = tf.keras.Input(shape=(14,22,1))

```
# convolutional part
#layer = tf.keras.layers.Conv2D(128, (4, 4), activation="relu")(input )
#layer = tf.keras.layers.Conv2D(64, (4, 4), activation="relu")(layer)
#layer = tf.keras.layers.Conv2D(32,(4,4), activation="relu")(layer)
#layer = tf.keras.layers.Conv2D(100, (5,13), activation="relu")(layer)
#layer = tf.keras.layers.Flatten()(layer)
```

```
# usual NN part remains the same
#layer = tf.keras.layers.Dense(50, activation="relu")(layer)
#layer = tf.keras.layers.Dense(4, activation="linear")(layer)
#model = tf.keras.Model(inputs=input,outputs=layer)
```

```
# continued on the next page ...
```

# before # now

Import torch instead of tensorflow

Just comment it out but do not delete yet (convenient to keep for a reference while we code the torch model ...)



## **PyTorch DAMPE example - model definition**

... edit model cnn dampe paper torch.py continued:

```
class DampeCNN(tr.nn.Module):
    def init (self):
        super(). init ()
        self.conv layer1 = tr.nn.Conv2d(1,
        self.conv_layer2 = tr.nn.Conv2d(128, 64, (4, 4))
        self.conv_layer3 = tr.nn.Conv2d(64, 32, (4,4))
self.conv_layer4 = tr.nn.Conv2d(32, 100, (5,13))
        self.flatten layer = tr.nn.Flatten()
        self.dense layer1 = tr.nn.Linear(100,50)
        self.dense layer2 = tr.nn.Linear(50, 4)
        self.relu
                            = tr.nn.ReLU()
    def forward(self, x):
        x = self.conv layer1(x)
        x = self.relu(x)
        x = self.conv layer2(x)
        x = self.relu(x)
        x = self.conv layer3(x)
        x = self.relu(x)
        x = self.conv layer4(x)
        x = self.relu(x)
        #print (x.shape)
        #raise SystemExit
        x = self.flatten layer(x)
        x = self.dense layer1(x)
        x = self.relu(x)
        x = self.dense layer2(x)
        return x
```

model = DampeCNN()

128, (4,4))

For simple models (including this one) we could in principle use a simplified syntax with tr.nn.Sequential (see our 2nd PyTorch installation test), but the low-level definition shown here offers more flexibility (for example in Transformer models that we will consider later)


### PyTorch DAMPE example - test

Let's test a bit our model. Similar to PyTorch case, we can feed it some random data and print the layer shape on the way. First, let's add a print statements to the model forward method in model cnn dampe paper torch.py:

```
class DampeCNN(tr.nn.Module):
    • • •
    • • •
    def forward(self, x):
        x = self.conv layer1(x)
        print ("Dimensions after layer 1:", x.shape)
        # etc
         • • •
```

Now let's run it the model a random data:

```
python -i model cnn dampe paper torch.py
>>> data=tr.tensor(np.zeros((100,1,14,22)),dtype=tr.float32)
>>> tmp = model(data)
Dimensions after layer 1: torch.Size([100, 128, 11, 19])
```





Result of a print statement



### PyTorch DAMPE example - test

### Now finally as we have all pieces in place, we can run the training of the DAMPE CNN model in PyTorch:

python run fit dampe torch.py

Processed valid: 99.9 elapsed:164.9s loss: 78.60382124119583 loss val: 21.724795972978747 Processed valid: 99.9 elapsed:165.9s loss: 16.86408121580472 loss\_val: 12.455400391741916 • • •

We will not run this training, just make sure that it works in principle ... Later you can run that CNN model in PyTorch and compare the performances with Tensorflow: execution time, convergence, final accuracy (loss) of the model ...

Now let's briefly re-cap our PyTroch CNN example ...





## Recap (PyTorch DAMPE example)

### run fit dampe torch.py

```
import torch as tr
import matplotlib.pyplot as plt
from model cnn dampe torch import model
from dampe import get dampe data
from pickle import dump
import time
data = get dampe data()
x, y = data['caloimages'], data['truthdata']
x, y = tr.tensor(x,dtype=tr.float32), tr.tensor(y,dtype=tr.float32)
x = x.squeeze(dim=3).unsqueeze(dim=1)
criterion = tr.nn.L1Loss()
optimizer = tr.optim.Adam(model.parameters())
# training loop
history = { 'loss':[], 'loss val':[] }
n total, n split, n batch = x.shape[0], int(x.shape[0] * 0.9), 32
for i in range(5): # loop over training epochs
    loss_train, loss_val, n_batches_train, n batches val, timestamp = 0, 0, 0, 0, time.time()
    model.train() # set model in training mode
    for j in range(0, n split, n batch): # loop over batches in training sub-sample
        print (f"Processed: {j*100./n split:3.1f}% elapsed:{time.time()-timestamp:4.1f}s",end="\r")
        optimizer.zero grad()
        y pred = model(x[j:j+n batch])
        loss = criterion(y_pred,y[j:j+n_batch])
        loss.backward()
        optimizer.step()
        loss train+=loss.item()
        n batches train+=1
    model.eval() # set model in validation mode
    for j in range(n split, n total, n batch): # loop over batches in validation sub-sample
        print (f"Processed valid: {(j-n split)*100./(n total-n split):3.1f} elapsed:{time.time()-timestamp:4.1f}s",end="\r")
        y pred = model(x[j:j+n batch])
        loss = criterion(y pred,y[j:j+n batch])
        loss val+=loss.item()
        n batches val+=1
    loss train, loss val = loss train / n batches train, loss val / n batches val
    print ("\nloss:",loss train, "loss val:",loss val)
    history['loss'], history['loss val'] = history['loss'] + [loss train], history['loss val'] + [loss val]
# save model and training history
tr.save(model.state dict(), './model cnn dampe paper torch.pth')
dump(history, open('history cnn dampe paper torch.p','wb')) # remains the same - we just changed the name of the file
```



## Recap (PyTorch DAMPE example)

### model\_cnn\_dampe\_paper\_torch.py

```
import torch as tr
class DampeCNN(tr.nn.Module):
   def init (self):
       super(). init ()
       self.conv layer1 = tr.nn.Conv2d(1, 128, (4,4))
       self.conv layer2 = tr.nn.Conv2d(128, 64, (4, 4))
       self.conv layer3 = tr.nn.Conv2d(64, 32, (4, 4))
       self.conv layer4 = tr.nn.Conv2d(32, 100, (5,13))
       self.flatten layer = tr.nn.Flatten()
       self.dense layer1 = tr.nn.Linear(100,50)
       self.dense layer2 = tr.nn.Linear(50, 4)
       self.relu = tr.nn.ReLU()
  def forward(self, x):
       x = self.conv layer1(x)
       x = self.relu(x)
       x = self.conv layer2(x)
       x = self.relu(x)
       x = self.conv layer3(x)
       x = self.relu(x)
       x = self.conv layer4(x)
       x = self.relu(x)
       x = self.flatten layer(x)
       x = self.dense_layer1(x)
       x = self.relu(x)
       x = self.dense layer2(x)
       return x
model = DampeCNN()
```



## Exercise (10 mins)

- Implement your previously developed DAMPE NN (not CNN) in PyTorch
  - → Goal: make sure you understand the layer implementation in PyTorch (where and how to look for documentation)
  - > You will create your own PyTorch model class similar to the CNN example
  - $\rightarrow$  No need to run the full training, just make sure that the training works and you are able to run  $\sim 1$  training epoch

From the previous example, you should already know the name of layers in PyTorch. For more information/references, see PyTorch documentation

• For example, the Flatten layer in PyTorch: <u>https://docs.pytorch.org/docs/stable/generated/torch.nn.Flatten.html</u>



### Part III: Transformers





### Transformers

Before starting preparing this tutorial, I promised myself not to make too obvious flat jokes connected to the comic books characters ...





### Transformers

Before starting preparing this tutorial, I promised myself not to make too obvious flat jokes connected to the comic books characters ...





February 14, 2023

*Now seriously... According to Stephen Wolfram:* 

### What Is ChatGPT Doing ... and Why Does It Work?



### It's Just Adding One Word at a Time

https://writings.stephenwolfram.com/2023/02/what-is-chatgpt-doing-and-why-does-it-work/





### On Large Language Models (LLMs)

From Wolfram's article:

The best thing about AI is its ability to	learn
	predict
	make
	understar
	do

And the remarkable thing is that when ChatGPT does something like write an essay what it's essentially doing is just asking over and over again "given the text so far, what should the next word be?"—and each time adding a word. (More precisely, as I'll explain, it's adding a "token", which could be just a part of a word, which is why it can sometimes "make up new words".)

Adding from myself: in some sense, LLMs can be considered ~ autocompletion on steroids



	4.5%
	3.5%
	3.2%
nd	3.1%
	2.9%



## On Large Language Models (LLMs) & Transformers

### **Computer Science > Computation and Language**

[Submitted on 12 Jun 2017 (v1), last revised 2 Aug 2023 (this version, v7)] **Attention Is All You Need** 

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin



It turned to be that the so-called "Transformer"-like architectures, originally suggested in the "Attention is All You Need" paper, are currently the best ones for LLM implementation (including GPT, Gemini, DeepSeek)

While the first GPT models were predicting tokens based on the preceding sequence of text of ~1-2 pages at most, GPT-4 turbo has a maximum sequence length (based on which it does "autocompletion") of the size of ~ Orwell's "1984" (or "The Hobbit", if you like)

https://arxiv.org/abs/1706.03762





### What is a token?

LLMs are prediction tokens one-by-one

- A word is token ? in most cases YES
- **Is token a word? NOT** necessarily  $\bullet$ 
  - A token can be a part of a word
  - →It can be a single digit, i.e. "1", "2", etc, combination of digits "12", "567", etc.

  - Event frequent combination of words, or entire sentences can be tokens, depending on the implementation



### The 72 tokens

This tool allows you to visualize the token s of a text prompt or token ization models of the various Google Cloud Ver tex AI Pa LM are also count e d , and hover ing over them will indicate their interr code of this application is available on Gi t hub .

→One single letter is a token, frequent combination of letters/syllables are tokens (that's how LLMs can invent new words)







### What is a token?

Neural Networks deal with numbers. How do we represent tokens as numbers? Say, GPT-4 uses ~100'000 token dictionary • Obviously, we assign every token a number from 0 to 100'000

- We can also represent those in vectors of 100'000 digits, all of which are 0 except for one:

[0,0,0, ..., 1, ..., 0, 0,0] \_\_\_\_\_\_ Token 1 ("Hello") 100'000 vectors [0,1,0, ..., 0, ..., 0, 0,0] **Token 2 ("world")** • • • [0,0,0, ..., 0, ..., 0, 0,1] — Token 100'000 ("bye")

100'000 numbers

Now what LLM does when prediction a next token, it essentially yields a so-called **logit** vector of 100'000 numbers each representing a probability for a certain token to be the next one in a sequence:

**# Numbers sum up to 1**  $[0.01, 0., 0.8, \ldots, 0.02, \ldots, 0, 0, 1]$ 

100'000 numbers

etc.



### Tokens & Embedding vectors

### **Embedding vectors (or Embeddings) - key concept in LLMs!**

- space (n is the dictionary size 100'000 in previous example) The answer is we don't do that!
- Instead, we **convert token into embedding vector** in lower dimension space through linear transformation:



• In previous illustration of token vectors - how do we encode similarity of tokens in that kind of vectors in that n-dimensional

embeddings will likely occur very aligned in the embedding space



### Tokens & Embedding vectors

eagle				
de camel	ər		duck	
t	bear	COW	chicken	
elephan	t	bird		
	dog	)	fish	
chimpanzee cheetah dolph	ca in	t		
		fly		
alligat	tor ile			
		ant		
		bee		

cranberry grape blueberry strawaspherry

apricot peach fig

> pineappleut cherry

> > mango

banana papaya melon apple

ribes avocado blackberry turnip

Projection of typical Embedding space in 2D (credit: Stephen Wolfram)

You can think of embedding space as a way of mapping language into a multidimensional cartesian space, where close-by-semanticalmeaning words/phrases/ syllables/etc. will normally appear as nearby vectors



### Data flow in Langauge Models



### **Embeddings-to-token** probabilities











## Tokens & Embedding vectors: try yourself

Before we go to Transformer implementation, let's first experiment with tokens & embeddings Why this is important: if we want to transform the use "Transformers" into physics applications we shall first get a clear idea of the "Transformers" original intended use (which is in fact language processing)

```
conda activate tr
conda install transformers
python
>>>
>>> from transformers import AutoTokenizer
>>> tokenizer = AutoTokenizer.from pretrained("gpt2")
>>> input tokens = tokenizer("I am GPT", return tensors="pt").input ids
>>> output tokens = tokenizer("Sono GPT", return tensors="pt").input ids
>>> print(input tokens)
  tensor([[ 40, 716, 402, 11571]])
>>> print(output tokens)
 tensor([[ 50, 29941, 402, 11571]])
>>> print (tokenizer.vocab size)
  50257
```

Don't close the python session yet, we continue on the next slide ...

Don't worry - we will implement tranformers ourself, we just need this package for token generation (we can/will do it ourself, but better to use the pre-cooked one)



- We will consider a **vastly oversimplified** example of language translation when one token encoding an English word is sought/trained to connect to exactly one other token that encodes the corresponding Italian word
- In this example, both English and Italian phrases are represented by 4 tokens, each token as a number is from 1 to 50256, 0 is reserved for empty token (masked in model predictions - as you will see further...)
- As can be seen the word "GPT" is represented presented by 2 tokens (402 and 11571)





### Tokens & Embedding vectors: try yourself

... continued from the previous slide



\* Tensor operations are supposed to be done in batches (e.g. batches of token sequences) for efficient parallel execution on the hardware (e.g. GPU). In our case we have only one sequence of four tokens, which explains the first dimension

We choose embedding dimension to be 512 - as in the original "Attention is all you need" paper (GPT-1 had 768)

Question: why do we need a dedicated PyTorch Embedding layer and why Linear is not OK for this job?

**Answer**: this is all about data format, Linear would fit the purpose if each toked would be presented in the format of [0,0,0, ...1, 0,0] *sparse vectors (of dimension vocab size) which is clearly not practical* ...









## A trivial Language Model: training

### Create a new file fit languate model eng it.py and run it:

```
import torch as tr
from transformers import AutoTokenizer
# tokenize data
tokenizer = AutoTokenizer.from pretrained("gpt2")
x = tokenizer("I am GPT", return tensors="pt").input ids
y = tokenizer("Sono GPT", return tensors="pt").input ids
# define the model, loss and optimizer for training
model = tr.nn.Sequential(
  tr.nn.Embedding(tokenizer.vocab size, 512),
  tr.nn.Linear(512,tokenizer.vocab size))
criterion = tr.nn.CrossEntropyLoss(ignore index=0)
optimizer = tr.optim.Adam(model.parameters())
# fit the model
model.train() # training mode on
for i in range(30):
    y \text{ pred} = \text{model}(x)
    loss = criterion(y pred.view(-1, tokenizer.vocab size), y.view(-1))
    # to understand the above view(...) transformation - uncomment the below:
    # print (y pred.shape, y.shape)
    # print (y pred.view(-1, tokenizer.vocab size).shape, y.view(-1).shape)
    loss.backward()
    optimizer.step()
    print ("Loss:", loss.item())
tr.save(model.state_dict(), 'model_eng_it.pth')
```



## A trivial Language Model: training

### Create a new file fit languate model eng it.py and run it:

```
import torch as tr
from transformers import AutoTokenizer
# tokenize data
tokenizer = AutoTokenizer.from pretrained("gpt2")
x = tokenizer("I am GPT", return tensors="pt").input ids
y = tokenizer("Sono GPT", return tensors="pt").input ids
# define the model, loss and optimizer for training
model = tr.nn.Sequential(
  tr.nn.Embedding(tokenizer.vocab size, 512),
  tr.nn.Linear(512, tokenizer.vocab size))
criterion = tr.nn.CrossEntropyLoss(ignore index=0)
optimizer = tr.optim.Adam(model.parameters())
# fit the model
model.train() # training mode on
for i in range(30):
    y \text{ pred} = \text{model}(x)
    loss = criterion(y pred.view(-1, tokenizer.vocab size), y.view(-1))
    # to understand the above view(...) transformation - uncomment the below:
    # print (y pred.shape, y.shape)
    # print (y pred.view(-1, tokenizer.vocab size).shape, y.view(-1).shape)
    loss.backward()
    optimizer.step()
    print ("Loss:", loss.item())
tr.save(model.state_dict(), 'model eng it.pth')
```



\*Note that CrossEntropyLoss of PyTorch automatically acts as SparseCategoricalCrossentropy of Tensorflow





### A trivial Language Model: inference

### Create a new file eval languate model eng\_it.py

```
import torch as tr
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from_pretrained("gpt2") # tokenizer instance
model = tr.nn.Sequential(
  tr.nn.Embedding(tokenizer.vocab size, 512),
  tr.nn.Linear(512, tokenizer.vocab_size)
model.load_state_dict(tr.load('model_eng_it.pth', weights_only=True)) # load trained model weights
model.eval() # set model in the evaluation mode
```

Now let's run it: python -i eval languate model eng it.py

```
>>> in_data = tokenizer("I am", return_tensors="pt").input_ids
>>> out data = model(in data)
>>> print (out data.shape)
>>> print (out_data)
  torch.Size([1, 2, 50257])
  tensor([[[-6.3526, -5.8213, -5.9804, ..., -6.5088, -5.9832, -6.1094],[-6.8054, -6.9886, -5.6512, ...,
  -6.5140, -6.8379, -6.2876]]],grad fn=<ViewBackward0>)
```

In principle we could run the Softmax transformation to convert the vector of dimension 50257 into the actual probabilities, however it is not needed since we are basically interested in the order only (Softmax is monotonic function) - at which position is the highest value?

- # define the model



## A trivial Language Model: inference

### Create a new file eval languate model eng it.py

```
import torch as tr
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from_pretrained("gpt2") # tokenize data
model = tr.nn.Sequential(
  tr.nn.Embedding(tokenizer.vocab size, 512),
  tr.nn.Linear(512, tokenizer.vocab size)
model.load_state_dict(tr.load('model_eng_it.pth', weights_only=True)) # load trained model weights
model.eval() # set model in the evaluation mode
```

Now let's run it: python -i eval languate model eng it.py

```
>>> in_data = tokenizer("I am", return_tensors="pt").input ids
>>> out data = model(in data)
>>> print (out data.shape)
>>> print (out_data)
  torch.Size([1, 2, 50257])
  -6.5140, -6.8379, -6.2876]]],grad fn=<ViewBackward0>)
>>> out indexmax = tr.argmax(out data, dim=-1)
>>> print (out_indexmax.shape)
  torch.Size([1, 2])
>>> print (out indexmax)
  tensor([[ 50, 29941]])
>>> print (tokenizer.batch decode(out indexmax))
   ['Sono']
```

# define the model

tensor([[[-6.3526, -5.8213, -5.9804, ..., -6.5088, -5.9832, -6.1094],[-6.8054, -6.9886, -5.6512, ...,

Et voilà! We see our Eng-It model in action!



## A trivial Language Model: inference (recap)

### Update eval languate model eng it.py to have everything in one place (will need it for later...)

```
import torch as tr
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from pretrained("gpt2") # tokenize data
model = tr.nn.Sequential(
                                               # define the model
  tr.nn.Embedding(tokenizer.vocab size, 512),
  tr.nn.Linear(512, tokenizer.vocab size)
model.load state dict(tr.load('model eng it.pth', weights only=True)) # load trained model
weights
model.eval() # set model in the evaluation mode
# example of model inference
in data = tokenizer("I am", return tensors="pt").input ids
out data = model(in data)
out indexmax = tr.argmax(out data, dim=-1)
human readable output = tokenizer.batch decode(out indexmax)
```



### Coming back to Transformers...





Now as we (hopefully) have some intuition regarding the Language models, we can get back to the Transformers. **Transformer is a model architecture.** In the previously considered example we used a trivial model architecture with 2 layers only:

- Embedding (tokens --> embedding vectors)
- Linear (embedding vectors --> token probabilities) ullet

```
model = tr.nn.Sequential(
  tr.nn.Embedding(tokenizer.vocab_size, 512),
  tr.nn.Linear(512, tokenizer.vocab_size)
```



### Coming back to Transformers...





Now as we (hopefully) have some intuition regarding the Language models, we can get back to the Transformers. **Transformer is a model architecture.** In the previously considered example we used a trivial model architecture with 2 layers only:

Embedding (tokens --> embedding vectors)



- Transformer
- Linear (embedding vectors --> token probabilities) ullet











Q = query, K = key

Self-attention



See also: https://fall-2023-python-programming-for-data-science.readthedocs.io/en/latest/Lectures/ Theme\_3-Model\_Engineering/Lecture\_20-Transformer\_Networks/Lecture\_20-Transformer\_Networks.html





Cross-attention



Q = query, K = key



See also: https://fall-2023-python-programming-for-data-science.readthedocs.io/en/latest/Lectures/ Theme\_3-Model\_Engineering/Lecture\_20-Transformer\_Networks/Lecture\_20-Transformer\_Networks.html







Both cross-attention and self-attention is used in Language translation problems (original model in "Attention is all you need" paper)

Self-attention











In chatbot models like GPT only self-attention is used (unless it deals with translation)





For clarity, let's not consider cross-attention and focus on self-attention only (all basic concepts are valid ~ identical for both cases)



In chatbot models like GPT only self-attention is used (unless you ask it for translation of something)





Input sequence of embeddings (e.g. "how", "are", "glacier", "caves", "formed") is transformed through attention matrix (attention block):



**X** - input; for now, think of **Q** = **K** = **V** = **X** (sequence of input embedding vectors)



Input sequence of embeddings (e.g. "how", "are", "glacier", "caves", "formed") is transformed through attention matrix (attention block):



**X** - input; for now, think of  $\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}$  (sequence of input embedding vectors)

For clarity, imagine that the Embedding space has only 1 dimension (not 512 or whatever...)





Input sequence of embeddings (e.g. "how", "are", "glacier", "caves", "formed") is transformed through attention matrix (attention block):









Input sequence of embeddings (e.g. "how", "are", "glacier", "caves", "formed") is transformed through attention matrix (attention block):

Output of the Transformer - same structure as the input: a sequence of embedding-space vectors (5 vectors in our example) but...

Each "embedding" is not merely the one corresponding to the word/token itself, but is rather a sum of all original embeddings weighted according to their correlation/attention







Input sequence of embeddings (e.g. "how", "are", "glacier", "caves", "formed") is transformed through attention matrix (attention block):

Output of the Transformer - same structure as the input: a sequence





## Attention: queries (Q), keys (K), values (V)

### Now where is the catch?

- Before we considered **Q=K=V** = (sequence of embedding vectors each one corresponding to an input token)

# $\mathbf{Q} = \mathbf{W}_{\mathbf{Q}} \cdot \mathbf{X} \quad \mathbf{K} = \mathbf{W}_{\mathbf{K}} \cdot \mathbf{X} \quad \mathbf{V} = \mathbf{W}_{\mathbf{V}} \cdot \mathbf{X}$

- - **Q** query-transformed "embeddings" (5 x 512 matrix ...)
  - **K** key-transformed "embeddings" (5 x 512 matrix ...)
  - V value-transformed "embeddings" (5 x 512 matrix ...)

If that would be the case, the only way were model could learn something is in the Embedding layer (the one that translates tokens into embedding vectors) - which would barely be enough for any real-life application

### • In reality is that Q, K, V are not merely the input embeddings, but linearly-transformed input embeddings!



W<sub>q</sub>, W<sub>k</sub>, W<sub>v</sub> - n x n matrices where n is the number of dimensions in embedding space (512 x 512 in our example) X - input embeddings (5 x 512 matrix in our example of "how", "are", "glacier", "caves", "formed")



## Attention: queries (Q), keys (K), values (V)

### Now where is the catch?

- Before we considered **Q=K=V** = (sequence of embedding vectors each one corresponding to an input token)

# $\mathbf{Q} = \mathbf{W}_{\mathbf{Q}} \cdot \mathbf{X}$ $\mathbf{K} = \mathbf{W}_{\mathbf{K}} \cdot \mathbf{X}$ $\mathbf{V} = \mathbf{W}_{\mathbf{V}} \cdot \mathbf{X}$

If that would be the case, the only way were model could learn something is in the Embedding layer (the one that translates tokens into embedding vectors) - which would barely be enough for any real-life application

### • In reality is that Q, K, V are not merely the input embeddings, but linearly-transformed input embeddings!



In PyTorch (given our example) the three matrices are implemented with the Linear layers:

- WQ = tr.nn.Linear(512, 512)
- W K = tr.nn.Linear(512, 512)
- W V = tr.nn.Linear(512, 512)


# Attention: queries (Q), keys (K), values (V)



• Okay, the transformer takes as an input a sequence of embedding corresponding to tokens ("how", "are", "glacier", "caves", "formed") • What do we want to predict? The answer is: for every input token we train the model to predict next one in the sequence...





# Transformer training: input & output



Remember: number if input and output tokens in self-attention models is always the same!

• Okay, the transformer takes as an input a sequence of embedding corresponding to tokens ("how", "are", "glacier", "caves", "formed") • What do we want to predict? The answer is: for every input token we train the model to predict next one in the sequence...





## Transformer training: input & output

• To avoid the model "looking into the future" at training, we mask elements below the main diagonal in the attention matrix





# Transformer training: input & output

• To avoid the model "looking into the future" at training, we mask elements below the main diagonal in the attention matrix



For example, the output of transformer for the word "glacier" will be a weighted sum of embeddings of word "glacier" itself and the words preceding in the sequence ("are","how") but not those after ("caves"). Otherwise, the model will just pick up (learn) to use the features of a word/token that comes right after - it will trivially work to predict n-1 tokens out of n, but will yield nonsense for the last one (as it has no information about the word that come after)





## Transformer inference: input & output

During the inference, we ask model to iteratively predict one toke at at time:

•	First call of the mod → Inputs: ["how"] → Outputs: ["are"]	del:		
•	Second call of the mo → Inputs: ["how", ' → Outputs: ["are", '	odel: " <b>are" ]</b> "glacier" ]		
•	Third call of the mod → Inputs: ["how", ' → Outputs: ["are", '	del: " <b>are",</b> "glacier",	" <b>glacier</b> " "caves"	<b>]</b>
•	Fourth call of the mo → Inputs: ["how", ' → Outputs: ["are", '	odel: " <b>are",</b> "glacier",	"glacier", "caves",	"cav "form

etc. (we can, in principle, continue forever ...)

res" ] med" ]

#### 113

# Transformer inference: input & output (cross-attention)

During the inference, we ask model to iteratively predict one toke at at time:

•	<pre>First call of the model:   → Inputs: ["how", "are", "glacier", "caves",   → Outputs: ["come"]</pre>
•	<pre>Second call of the model:   → Inputs: ["how", "are", "glacier", "caves",   → Outputs: ["come","si"]</pre>
•	<pre>Third call of the model:   → Inputs: ["how", "are", "glacier", "caves",   → Outputs: ["come","si", "glaciali"]</pre>
•	<pre>Fourth call of the model:   → Inputs: ["how", "are", "glacier", "caves",   → Outputs: ["come","si", "glaciali", "grotte"</pre>
•	<pre>Fifth call of the model (we don't predict, we t → Inputs: ["how", "are", "glacier", "caves", → Outputs: ["come","si", "glaciali", "grotte"</pre>
•	<pre>Fifth call of the model (we don't predict, we t → Inputs: ["how", "are", "glacier", "caves", → Outputs: ["come","si", "glaciali", "grotte"</pre>

If trained well, the model will likely yield a stop-sequence token telling us that it no more tokens are needed and the English sentence translation is completed









### Enough (almost) theory, let's get back to coding...

ion at the end sadd back the desel

reod.use is a false

and suscerve structure

od.use z e False

eod\_use\_x = False

eod.use\_y = False

mod\_use\_z = True

ERATOR CLASSES -

select-1

don see THIRKOR Zhi

["please select exactly two objects, "

Operator): (propiet of the selected object\*\*\*

### Attention block code

#### Create transformer torch.py :

```
import torch as tr
import math
class Attention(tr.nn.Module):
    def init (self, d model):
        super(). init ()
       # Linear layers for transforming inputs
        self.W q = tr.nn.Linear(d_model, d_model) # Query transformation
        self.W k = tr.nn.Linear(d model, d model) # Key transformation
        self.W v = tr.nn.Linear(d model, d model) # Value transformation
        self.d model = d model # embedding dimension
    def scaled dot product attention(self, Q, K, V, mask=None):
        # Calculate attention scores
        attn scores = tr.matmul(Q, K.transpose(-2, -1)) / math.sqrt(self.d model)
        # Apply mask (for example the under-diagonal elements)
       if mask is not None: attn scores = attn scores.masked fill(mask==0,-1e9)
        # Softmax is applied to obtain attention probabilities
        attn probs = tr.softmax(attn scores, dim=-1)
        # Multiply by values to obtain the final output
        return tr.matmul(attn probs, V)
    def forward(self, x, mask=None):
        # Apply linear transformations and split head
       Q = self.W q(x)
        K = self.W k(x)
       V = self.W v(x)
        # Perform scaled dot-product attention and return
        attn output = self.scaled dot product attention(Q, K, V, mask)
       return attn output
```

We already know W\_k, W\_q, W\_v matrices

- Calculate attention (QK) matrix
- Apply mask (for example QK elements below the main diagoanal)
- Obtained QK matrix embedding vector (divided by d model to normalize typical value and avoid growth with large d model)
- Softamax is applied to normalize to QK matrix element to probabilities

Propagate our input sequence through the attention block











### Attention block code: test

#### Create test attention block torch.py :

```
import torch as tr
from transformer torch import Attention
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from_pretrained("gpt2")
attentionblock = Attention(512)
x = tr.randint(1, tokenizer.vocab_size-1, (32, 100))
print (x.shape)
x embed = tr.nn.Embedding(tokenizer.vocab size, 512)(x)
print (x_embed.shape)
x_attention = attentionblock(x_embed)
```

```
print (x embed.shape)
```

python test\_attention\_block\_torch.py :

```
torch.Size([32, 100])
torch.Size([32, 100, 512])
torch.Size([32, 100, 512])
```

We generate a batch of random data 32 sequences of 100 tokens each. Then we convert it to embedding vectors

We process the embedding vectors of our input data through the attention block, in the end we should get the output sequence if the same dimension

### 2



# Exercise (5 mins): understanding of the Attention

#### Look inside the function scaled dot product attention of the Attention class

- → Insert print statements to show the dimension (shape) of Q, K, V, attn scores, attn probs ... Do you understand those?
- Print the values of attn scores and attn probs matrices
  - ... Does the second one look like probabilities?

#### 118

• In reality, each transformer block has multiple attention heads. What does that mean? Consider our example:





• In reality, each transformer block has multiple attention heads. What does that mean? Consider our example:







Vectors in embedding space are split into equal size parts processed separately each by its own Transformer (Key-Query-Value transformation), which is called a "head". In a canonical "Attention is all you need" paper, there are 8 heads per attention block, which makes the dimension of sub-embedding for each head 512/8 = 64.





Vectors in embedding space are split into equal size parts processed separately each by its own Transformer (Key-Query-Value transformation), which is called a "head". In a canonical "Attention is all you need" paper, there are 8 heads per attention block, which makes the dimension of sub-embedding for each head 512/8 = 64.

In practice, the algorithm is as follows:

- 1. Apply key/value/query transformation to each embedding
- 2. Split embeddings into multiple equal-length parts (8 pars, 64 each in the original paper)
- 3. Apply attention transformation for each sub-embedding (head) (8 attention blocks)

4. Combine the outputs of 8 transformations back to the original-size embedding vector (dimension 512 in our example)



Main advantage if multi-head attention is that at the ~same computation complexity, the model can encode different attention patterns in through different heads in one attention block



For example, the first head may converge towards focusing on correlation in the beginning and the end of the phrase, putting more weight to those ("how" and "caves" in our example), while the other head may focus more on attention patters between subsequent words ("are", "glacier") etc. Remember that these attention patterns are encoded in the W<sub>Q</sub>, W<sub>K</sub>, W<sub>V</sub>, matrices...



## Multi-head attention block: code

#### Now let's update the transformer torch.py code (make a back-up copy of it before!):

```
#class Attention(tr.nn.Module):
                              # before
class MultiHeadAttention(tr.nn.Module): # now
                               # before
   #def init (self, d model):
   def init (self, d model, num heads): # now
       super(). init ()
       ... everything before remains the same, add the following in the end of _____init___:
       self.num heads = num heads
       self.W o = tr.nn.Linear(d model, d model) # Mixing matrix (we'll understand it later)
   def scaled dot product attention(self, Q, K, V, mask=None):
       ... everything the same except for replacing "d model" with "d k":
       #attn scores = tr.matmul(Q, K.transpose(-2, -1)) / math.sqrt(self.d model)
       attn scores = tr.matmul(Q, K.transpose(-2, -1)) / math.sqrt(self.d k)
   def forward(self, x, mask=None):
       # add the following code to split heads after the Q,V,K definition
       Q = self.split heads(self.W q(Q))
       K = self.split heads(self.W k(K))
       V = self.split_heads(self.W_v(V))
       # add the following after the "attn output" definition and before "return":
       attn output = self.combine heads(attn output)  # combine heads
       attn output = self.W o(attn output)
```

... contiunued on the next pate (we need to implement "split\_heads" and "combine\_heads" methods )

# Number of attention heads

- # mix heads

We need to define number of heads and also to create a "mixing" matrix to mix the heads in the end

Replace "d\_model" with "d\_k"

We first split the Q, K, V vectors into parts (8 parts 64 length each in our example) to be processed by 8 independent Attention transformations

Then we combine back the output





## Multi-head attention block: code

#### edit transformer torch.py (continued from the previous slide):

```
# ... continued from the previous page
# add these methods :
def split heads(self, x):
    # Reshape the input to have num heads for multi-head attention
    batch size, seq length, d model = x.shape
    print ("Tensor shape before splitting heads", x.shape)
    print ("Tensor shape after splitting heads", output.shape)
    return output
```

```
def combine heads(self, x):
```

```
# Combine the multiple heads back to original shape
batch size, , seq length, d k = x.shape
output = x.transpose(1, 2).contiguous()
output = output.view(batch size, seq length, self.d model)
print ("Tensor shape before combining heads", x.shape)
print ("Tensor shape after combining heads", output.shape)
return output
```





## Multi-head attention block: code (recap and test)

#### transformer\_torch.py:

```
import torch as tr
import math
class MultiHeadAttention(tr.nn.Module):
   def init (self, d model, num heads):
       super(). init ()
       # Linear layers for transforming inputs
       self.W q = tr.nn.Linear(d model, d model) # Query transformation
       self.W k = tr.nn.Linear(d model, d model) # Key transformation
       self.W v = tr.nn.Linear(d model, d model) # Value transformation
       self.d model = d model
                                           # embedding dimension
       self.num heads = num heads
                                         # Number of attention heads
       self.d k = d model // num heads  # Dimension of each head's key, query, and valu
       self.W o = tr.nn.Linear(d model, d model)# Mixing matrix (we'll understand it later)
   def scaled dot product attention(self, Q, K, V, mask=None):
       # Calculate attention scores
       attn scores = tr.matmul(Q, K.transpose(-2, -1)) / math.sqrt(self.d k)
       # Apply mask (for example the under-diagonal elements)
       if mask is not None: attn scores = attn scores.masked fill(mask==0,-1e9)
       # Softmax is applied to obtain attention probabilities
       attn probs = tr.softmax(attn scores, dim=-1)
       # Multiply by values to obtain the final output
       return tr.matmul(attn probs, V)
   def forward(self, x, mask=None):
       # Apply linear transformations and split head
       Q = self.W q(x)
       K = self.W k(x)
       V = self.W v(x)
       Q = self.split heads(self.W q(Q))
       K = self.split heads(self.W k(K))
       V = self.split heads(self.W v(V))
       # Perform scaled dot-product attention and return
       attn output = self.scaled dot product attention(Q, K, V, mask)
       attn output = self.combine heads(attn output) # combine heads
       attn output = self.W o(attn output)
                                                      # mix heads
       return attn output
   def split heads(self, x):
       # Reshape the input to have num heads for multi-head attention
       batch size, seq length, d model = x.shape
       output = x.view(batch size, seq length, self.num heads, self.d k).transpose(1, 2)
       print ("Tensor shape before splitting heads", x.shape)
                                                                     # remove me later!
       print ("Tensor shape after splitting heads", output.shape)
                                                                     # remove me later!
       return output
   def combine heads(self, x):
       # Combine the multiple heads back to original shape
       batch_size, _, seq_length, d k = x.shape
       output = x.transpose(1, 2).contiguous()
       output = output.view(batch_size, seq_length, self.d_model)
                                                                    # remove me later!
       print ("Tensor shape before combining heads", x.shape)
       print ("Tensor shape after combining heads", output.shape)
                                                                    # remove me later!
       return output
```

#### In test\_attention\_block\_torch.py replace:

```
#from transformer_torch import Attention
from transformer_torch import MultiHeadAttention
....
#attentionblock = Attention(512)
attentionblock = MultiHeadAttention(512, 8)
....
```

#### Run the test:





### Multi-head attention block: code recap and test

#### transformer torch.py:

```
import torch as tr
import math
class MultiHeadAttention(tr.nn.Module):
   def init (self, d model, num heads):
       super(). init ()
       # Linear layers for transforming inputs
       self.W q = tr.nn.Linear(d model, d model) # Query transformation
       self.W k = tr.nn.Linear(d model, d model) # Key transformation
       self.W_v = tr.nn.Linear(d_model, d_model) # Value transformation
       self.d model = d model
                                            # embedding dimension
       self.num heads = num heads
                                          # Number of attention heads
       self.d k = d model // num heads  # Dimension of each head's key, query, and valu
       self.W o = tr.nn.Linear(d model, d model)# Mixing matrix (we'll understand it later)
   def scaled dot product attention(self, Q, K, V, mask=None):
       # Calculate attention scores
       attn_scores = tr.matmul(Q, K.transpose(-2, -1)) / math.sqrt(self.d_k)
       # Apply mask (for example the under-diagonal elements)
       if mask is not None: attn scores = attn scores.masked fill(mask==0,-1e9)
       # Softmax is applied to obtain attention probabilities
       attn probs = tr.softmax(attn scores, dim=-1)
       # Multiply by values to obtain the final output
       return tr.matmul(attn probs, V)
   def forward(self, x, mask=None):
       # Apply linear transformations and split head
       Q = self.W q(x)
       K = self.W k(x)
       V = self.W v(x)
       Q = self.split heads(self.W q(Q))
       K = self.split heads(self.W k(K))
       V = self.split heads(self.W v(V))
       # Perform scaled dot-product attention and return
       attn output = self.scaled dot product attention(Q, K, V, mask)
       attn output = self.combine heads(attn output) # combine heads
       attn output = self.W o(attn output)
                                                      # mix heads
       return attn output
   def split heads(self, x):
       # Reshape the input to have num heads for multi-head attention
       batch size, seq length, d model = x.shape
       output = x.view(batch size, seq length, self.num heads, self.d k).transpose(1, 2)
       print ("Tensor shape before splitting heads", x.shape)
                                                                     # remove me later!
       print ("Tensor shape after splitting heads", output.shape)
                                                                     # remove me later!
       return output
   def combine heads(self, x):
       # Combine the multiple heads back to original shape
       batch size, , seq length, d k = x.shape
       output = x.transpose(1, 2).contiguous()
       output = output.view(batch size, seq length, self.d model)
       print ("Tensor shape before combining heads", x.shape)
                                                                    # remove me later!
                                                                    # remove me later!
       print ("Tensor shape after combining heads", output.shape)
       return output
```

**MultiHeadAttention** is a building block of a Transformer model. Now let's build the model itself - which is, roughly speaking, just a sequence of multiple instances of MultiHeadAttention stacked one after another (sort of convolutional filters in a CNN)...



Original cross-attention Transformer from the "Attention is all you need " paper <a href="https://arxiv.org/abs/1706.03762">https://arxiv.org/abs/1706.03762</a>



See the full code of cross-attention tutorial in: <u>https://</u> www.datacamp.com/tutorial/building-a-transformer-with-py-torch





**Self-attention-only** (decoder-only) Transformer used in GPT-like models



### So-called "Decoder-only" architecture is Considered in this tutorial





**Self-attention-only** (decoder-only) Transformer used in GPT-like models

**Transformer Layer** 

MultiHeadAttention **block** that we have already implemented in PyTorch





### **Self-attention-only** (decoder-only) Transformer used in GPT-like models





### Transformer layer: code

#### Edit transformer torch.py (add the following in the end):

```
class TransformerLayer(tr.nn.Module): # a.k.a. "Decoder" Layer
   def init (self, d model, num heads, d ff, dropout):
       super(). init ()
       # instance of attnetion head
       self.self attn = MultiHeadAttention(d model, num heads)
       # instance of layer normalisation
       self.norm1 = tr.nn.LayerNorm(d model)
       # layers of feed-forward part
       self.fc1 = tr.nn.Linear(d model, d ff)
       self.fc2 = tr.nn.Linear(d ff, d model)
       self.relu = tr.nn.ReLU()
       # another layer normalisation and dropout
       self.norm3 = tr.nn.LayerNorm(d model)
       self.dropout = tr.nn.Dropout(dropout)
   def forward(self, x, mask):
       # attention block
       attn output = self.self attn(x, mask)
       # sum attention output and the original input, normalise
       x = self.norm1(x + self.dropout(attn output))
       # non-linear feed-forward oparation
       ff output = self.fc2(self.relu(self.fc1(x)))
       # sum the non-linear output and the with its input, normalise
       x = self.norm3(x + self.dropout(ff output))
       return x
```





## Transformer layer: code

#### Edit transformer torch.py (add the following in the end):

```
class TransformerLayer(tr.nn.Module): # a.k.a. "Decoder" Layer
   def init (self, d model, num heads, d ff, dropout):
        super(). init ()
       # instance of attnetion head
       self.self attn = MultiHeadAttention(d model, num heads)
       # instance of layer normalisation
       self.norm1 = tr.nn.LayerNorm(d model)
       # layers of feed-forward part
       self.fc1 = tr.nn.Linear(d model, d ff)
       self.fc2 = tr.nn.Linear(d ff, d model)
       self.relu = tr.nn.ReLU()
        # another layer normalisation and dropout
       self.norm3 = tr.nn.LayerNorm(d model)
       self.dropout = tr.nn.Dropout(dropout)
   def forward(self, x, mask):
        # attention block
       attn output = self.self attn(x, mask)
       # sum attention output and the original input, normalise
       x = self.norm1(x + self.dropout(attn output))
       # non-linear feed-forward oparation
       ff output = self.fc2(self.relu(self.fc1(x)))
       # sum the non-linear output and the with its input, normalise
       x = self.norm3(x + self.dropout(ff output))
       return x
```





Pay attention to flow of non-changed input that is always added to the output (inspired by ResNET convolutional nets...)



## **Positional Encoding**

So far we never mentioned positional encoding... How does the Attention transformation know about relative positions of tokens/ embedding in the sequence? A neat trick to facilitate positional encoding is to add to the embedding vector another vector that encodes the position (Positional Encoding a.k.a. Positional Embedding). This is done before the data is fed to transformer layers.

Position encoding vector:

where  $PE(pos, i) = \begin{cases} sin(pos/1000^{i/d}), \text{ for } i=0,2,...\\ cos(pos/1000^{(i-1)/d}), \text{ for } i=1,3,... \end{cases}$ 

Why 10000?  $\rightarrow$  Wide range of wavelets covering patterns from short to long distances. For example, for small i the period is close to ~1 - short dependancies, for the largest i period gets very high - close to 10000 (long dependencies)









## Positional Encoding

Why this specific definition is used?



2. Values limited in +-1 range  $\rightarrow$  good for NN



Embedding vector index (*i*)

#### 3. Encodes relative distance\*:



\* Dot product between between two positional vectors is only defined by their relative position difference



# Positional Encoding: code implementation

#### Edit transformer torch.py (add the following in the end):

```
class PositionalEncoding(tr.nn.Module):
   def init (self, d model, max seq length):
       super(). init ()
```

```
pe = tr.zeros(max seq length, d_model)
position = tr.arange(0, max seq length, dtype=tr.float).unsqueeze(1)
div term = tr.exp(tr.arange(0, d model, 2).float() / d model * -(math.log(10000.0)))
```

```
pe[:, 0::2] = tr.sin(position * div term)
pe[:, 1::2] = tr.cos(position * div term)
```

```
# register pe as part of the model but do not trat as (trainable) parameters
self.register buffer('pe', pe)
```

```
def forward(self, x):
   return x + self.pe[:x.size(1)]
```

#### Let's test it, run python -i transformer torch.py:

```
>>> p = PositionalEncoding(512, 100) # embedding dimension 512, maximum 100 tokens in a sequence
>>> x = tr.zeros(100, 512) # 100 embedding-size (512) vectors with with initiated with zeros
>>> print (position embeddings[0])
  tensor([0., 1., 0., 1., 0., ... 1., 0., 1.])
>>> print (position embeddings[1])
  tensor([8.4147e-01, 5.4030e-01,..., 1.0366e-04, 1.0000e+00])
```

>>> position embeddings = p(x) # this is a sum of embedding itself + position, however our embeddings are just zeros...



# Putting Transformer pieces together

#### Incorporate everything into the final model $\rightarrow$ edit transformer torch.py (add the following in the end):

```
class Transformer(tr.nn.Module):
    def init (self, tgt vocab size, d model, num heads, num layers, d ff, max seq length, dropout):
        super(). init ()
        self.embedding = tr.nn.Embedding(tgt_vocab_size, d_model)
        self.positional encoding = PositionalEncoding(d model, max seq length)
        self.layers = tr.nn.ModuleList([TransformerLayer(d model, num heads, d ff, dropout) for in range(num layers)])
        self.fc = tr.nn.Linear(d_model, tgt_vocab_size) # conversion from embedding space to token probabilities
        self.dropout = tr.nn.Dropout(dropout)
    def generate mask(self, x): # mask 0 (empty) tokens and above-diagonal elements during the attention matrix generation
       mask nonzero = (x != 0).unsqueeze(1).unsqueeze(3)
        seq length = x.shape[1]
       mask = (1 - tr.triu(tr.ones(1, seq length, seq length), diagonal=1)).bool()
       mask = mask nonzero & mask
        return mask
    def forward(self, x):
       mask = self.generate mask(x)
        x embedded = self.dropout(self.positional encoding(self.embedding(x)))
        output = x embedded
        for layer in self.layers:
            output = layer(output, mask)
        return self.fc(output)
```



## Training Transformer with mockup data

#### Now let's test the code with a random data, create train transformer torch.py:

```
import torch as tr
from transformer torch import Transformer
tgt vocab size = 5000
d model = 512
num heads = 8
num layers = 6
d ff = 2048
max seq length = 100
dropout = 0.1
transformer = Transformer(tgt_vocab_size, d_model, num_heads, num_layers, d_ff, max_seq_length, dropout)
x = data[:, :-1]  # remember, we remove last token in the sequence
y = data[:, 1:] # the target for training is the same as the input but shifted by +1
criterion = tr.nn.CrossEntropyLoss(ignore index=0)
optimizer = tr.optim.Adam(transformer.parameters()) #, lr=0.0001, betas=(0.9, 0.98), eps=1e-9)
transformer.train()
for epoch in range(10):
    optimizer.zero grad()
    y pred = transformer(x)
    loss = criterion(y pred.contiguous().view(-1, tgt vocab size), y.contiguous().view(-1))
    loss.backward()
    optimizer.step()
    print(f"Epoch: {epoch+1}, Loss: {loss.item()}")
```

#### Run the test: python train transformer\_torch.py:

```
Epoch: 1, Loss: 8.696113586425781
Epoch: 2, Loss: 8.363970756530762
• • •
```

data = tr.randint(1, tgt\_vocab\_size, (64, max\_seq\_length)) # radom integers in the range from 1 to tgt\_vocab\_size (batch\_size, seq\_length)



# Training Transformer with real data

• Let's run it with some real text, create train transformer torch real.py:

```
import torch as tr
from transformer torch import Transformer
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from pretrained("gpt2")
tgt vocab size = tokenizer.vocab size
d model = 128 #512 # we use simpler model compared to the original "Attention is all you need" paper
num heads = 8
num layers = 2 #6 # same here .. simple model
d ff = 2048
max seq length = 10 \# 100 \# and here...
dropout = 0.1
transformer = Transformer(tgt vocab size, d model, num heads, num layers, d ff, max seq length, dropout)
with open('little prince.txt', 'r') as f: text = f.read()
data = tokenizer(text, return tensors="pt").input ids[0]
n seq = 32
data = data[:n_seq*max_seq_length] # use a small portion of text, 32*100 tokens
data = data.view(n seq,max seq length)
x = data[:,:-1]
y = data[:, 1:]
criterion = tr.nn.CrossEntropyLoss(ignore index=0)
optimizer = tr.optim.Adam(transformer.parameters())
transformer.train()
for epoch in range(50):
    optimizer.zero grad()
    y pred = transformer(x)
    loss = criterion(y pred.contiguous().view(-1, tgt vocab size), y.contiguous().view(-1))
    loss.backward()
    optimizer.step()
    print(f"Epoch: {epoch+1}, Loss: {loss.item()}")
    tr.save(transformer.state dict(), './model.pth')
```

- Run the training python train transformer torch real.py:

Epoch: 50, Loss: 0.7570830583572388

• Create little prince.txt manually by copying the first fragment of text (few pages) from https://archive.org/stream/TheLittlePrince-English/littleprince\_djvu.txt





### Transformer prediction/inference: word-by-word

• Create test transformer torch real.py:

```
import torch as tr
from transformer_model import Transformer
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from pretrained("gpt2")
tgt vocab size = tokenizer.vocab size
d model = 128
num heads = 8
num layers = 2
d ff = 2048
max seq length = 10
dropout = 0.1
transformer = Transformer(tgt vocab size, d model, num heads, num layers, d ff, max seq length, dropout)
transformer.load state dict(tr.load('./model.pth'))
transformer.eval()
def predict next word(x):
   print ("Input string: ", x)
   data = tokenizer(x, return tensors="pt").input ids
   output = transformer(data) # apply our transformer model
   outtokens = tr.argmax(output,dim=-1)
   return tokenizer.batch_decode(outtokens[:,-1:])[0] # decode only last token
```

• Run the code: python test transformer\_torch\_real.py:

```
>>> x="The Little Prince"
>>>
>>> x+=predict next word(x)
Input string: The Little Prince
>>> x+=predict next word(x)
Input string: The Little Prince appears
>>> x+=predict next word(x)
Input string: The Little Prince appears to
>>> x+=predict next word(x)
Input string: The Little Prince appears to be
>>> x+=predict next word(x)
Input string: The Little Prince appears to be a
• • •
```



# Transformers & Language Models: summary

- Language Model ~ a neural network that predicts word/token one at a time given the tokens before
- Embedding key instrument in Language Models:
   representation of words/tokens with continues vectors in multidimensional space
- Transformers- the most powerful/efficient architecture for language models so far
- Attention transformation: central component of Transformers
   alike convolutional filter in CNNs
- Can we use Transformers outside Language Models? Of course!





### Step 4: Visualizing the results

prediction of layer 9





AI generated



# **Final Exercise:** Transformer in physics (DAMPE case)

### Consider applying our Transformer model to the analysis of DAMPE data. The goal is to train the model to predict the development of particle interaction (shower) in the detector.

- Consider DAMPE calorimeter image (14 x 22) as a sequence of (14) "words"
- Can we predict next word based on the current and previous ones?

NOTE: for simplicity, we have specifically chosen the task that somewhat resembles the Language Model text sequence. However, Transformers can be equally applied to any other task in physics, such as classification (identifying particle type from the image), regression (predicting particle direction) etc.

DAMPE

>>> from dampe import get dampe data, plot dampe event >>> data = get dampe data() >>> plot dampe event(data['caloimages'][12])







# **Final Exercise:** Transformer in physics (DAMPE case)

### Consider applying our Transformer model to the analysis of DAMPE data. The goal is to train the model to predict the development of particle interaction (shower) in the detector.

- $\rightarrow$  Consider DAMPE calorimeter image (14 x 22) as a sequence of (14) "words"
- Can we predict next word based on the current and previous ones?

### **Implementation hints:**

- No need for tokenizer
- In the Transformer model use Linear layer instead of Embedding to convert 22-pixel array ("word") into 512 vector
- Output of the model should be a set of 22 numbers, each representing a signal in a pixel of subsequent layer (instead of probabilities)
- Use L1Loss instead of CrossEntropyLoss



>>> from dampe import get dampe data, plot dampe event >>> data = get dampe data() >>> plot dampe event(data['caloimages'][12])






## Step 1: Modify Transformer to DAMPE case

### Create a DAMPE transformer file: cp transformer torch.py transformer\_torch\_dampe.py, edit it:

class Transformer(tr.nn.Module): def init (self, tgt vocab size, d model, num heads, num layers, d ff, max seq length, dropout):

```
#self.embedding = tr.nn.Embedding(tgt vocab size, d model)
self.embedding = tr.nn.Linear(tgt vocab size, d model)
. . . .
```

```
def generate mask(self, x): #
    #mask nonzero = (x != 0).unsqueeze(1).unsqueeze(3)
    seq length = x.shape[1]
   mask = (1 - tr.triu(tr.ones(1, seq length, seq length), diagonal=1)).bool()
    #mask = mask nonzero & mask
    return mask
```

Replace embedding layer with with a linear transform, tgt vocab size will 22 in our case (one row of DAMPE image). We still call it "embedding" since the concept is the same -- we map our "word" (array of 22 numbers) into a vector in a higher-dimensional space

Comment our the mask nonzero part of the mask, since we will not have such thing as empty tokens. Also the input data format changed before we had one number that is a token ID, now we have an array of 22 numbers which represents a "word" (row of DAMPE pixels)







## Step 2: Modify torch training code to DAMPE case

### Create training code: cp run fit dampe torch.py run fit dampe torch transformer.py, editit:

```
#from model cnn dampe torch import model
from transformer torch dampe import Transformer
model = Transformer(tgt vocab size=22, d model=64, num heads=2, num layers=2, d ff=256, max seq length=13, dropout=0.)
#model = Transformer(tgt vocab size=22, d model=512, num heads=8, num layers=6, d ff=2048, max seq length=13, dropout=0.)
#x, y = data['caloimages'], data['truthdata']
x = data['caloimages'][:,:-1,:]
y = data['caloimages'][:,1:, :]
#x = x.squeeze(dim=3).unsqueeze(dim=1)
x = tr.squeeze(x, dim=3)
y = tr.squeeze(y, dim=3)
. . .
```

tr.save(model.state\_dict(), './model transformer dampe.pth') #'./model cnn dampe paper torch.pth') dump(history, open('history transformer dampe.p','wb')) #'history cnn dampe paper torch.p'

Initiate the transformer model. We will use simpler version of the model (compared to the original one), having 2 heads instead of 8, 2 layers instead of 6, d\_model = 64 instead of 512, and d\_ff=256 instead of 2048

Unlike before, our images are now both input (x) and target (y). The difference is that in the input we remove the last word (last row of image, so now it has 13 rows, not 14). For the output, we shift the rows by +1

We change the shape of the data to remove the last dimension/mode, so that instead of [n samples, n rows, n columns, 1] it becomes [n samples, n rows, n columns] Remember that the last axis was required for CNNs and represented the number of image channels (1 in our DAMPE case, 3 in case of typical RGB color images)

> Don't forget to rename the model and history to save ...







# Step 3: Running torch training code

### Set number of epochs to 10 in run fit dampe torch transformer.py and run the code:

pyhon run fit dampe torch transformer.py Processed valid: 99.9 elapsed:28.5s loss: 0.03037726884624309 loss val: 0.026196253926468058 Processed valid: 99.9 elapsed:28.8s loss: 0.022899649906208453 loss val: 0.019884693373397395 Processed valid: 99.9 elapsed:28.6s loss: 0.01885740954644424 loss val: 0.017663721486019926 Processed valid: 99.9 elapsed:28.5s loss: 0.016504979629425022 loss val: 0.015762685200300167 Processed valid: 99.9 elapsed:28.6s loss: 0.015107720140568331 loss val: 0.014291047355629973 Processed valid: 99.9 elapsed:28.6s loss: 0.01421325063580615 loss val: 0.013719815335043514 Processed valid: 99.9 elapsed:28.6s loss: 0.013527586701636979 loss val: 0.013062344417227683 Processed valid: 99.9 elapsed:28.3s loss: 0.012988961886007207 loss val: 0.01236261014600058 Processed valid: 99.9 elapsed:28.5s loss: 0.012554175169350597 loss val: 0.012217357914122913 Processed valid: 99.9 elapsed:28.2s loss: 0.012218388652761583 loss val: 0.011963540029519045

Depending on the hardware, it may take 5 to 10 minutes to finish. We have this time to discuss, ask question etc. If it takes significantly longer on your hardware, either use small number of iterations ( $\sim 2$ ) - it will be enough for illustrative purposes, or just copy the already pre-trained model from the DAMPE package you dowloaded earlier (see saved models folders there)





cp run fit dampe torch transformer.py test dampe torch transformer.py, edit:test dampe torch transformer.py:

```
\#criterion = tr.nn.L1Loss() \# < --- remove everything below this line (training loop etc.)
# ... add the following code instead:
from dampe import plot dampe event # we will need it for plotting
model.load state dict(tr.load("model transformer dampe.pth"))
model.eval() # set the model into evaluation mode
def predict dampe layers(x, event id, predict from layer):
   in image = x[event id:event id+1,:predict from layer] # DAMPE image truncated after [predict from layer]
    plot dampe event(in image[0],
                                                          # show image before any prediction
    title=f"input truncated image (no predictions yet")
    for i in range(predict from layer, 14):
                                                          # loop over remaining layers
        out image = model(in image)
                                                          # predicted image shifted by +1
        out_last_predicted_layer = out_image[:,-1:,:]
                                                          # take the last layer of the predicted image and ...
        in image = tr.cat(
                                                          # ... add it to the input image
            (in image, out last predicted layer), dim=1)
        plot_dampe_event(in_image.detach().numpy()[0],
                                                          # show image after one prediction step
        title=f"prediction of layer {i}")
    return in image.detach().numpy()[0]
# plotting code
event id = 12
                              # let's pick some event from the data
predict from layer = 9 # let's feed the model an image of 9 layers and see if it can predict the remaining 5
ai generated image = predict dampe layers(x, event id, predict from layer)
plot_dampe_event(data['caloimages'][event id],title="that is the original image")
plot dampe event(ai generated image,title=f"that image is generated starting from layer {predict from layer}")
```





input truncated image (no predictions so far)





prediction of layer 10





AI generated



prediction of layer 11





AI generated



prediction of layer 12





prediction of layer 13









if its >repo\_path = prop\_asting / for free\_path if its >repo\_path) if its its is its >repo\_path) if its its is its i











### matpletlib





