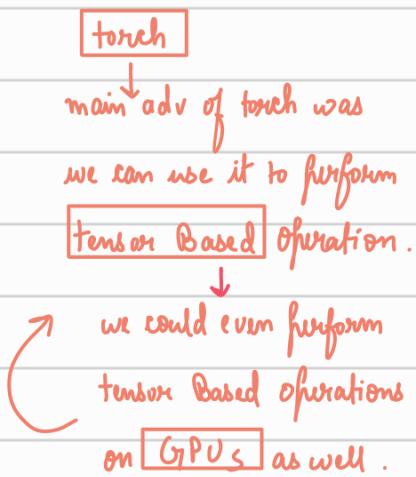


## Origin Story of Pytorch.

In [2002] → very powerful scientific tool  
came into picture



Although torch started as scientific computing framework but all researchers thought if we can apply tensor Based operations on GPUs it means we can use this framework torch to build all applications as well.

Then they wrote many NN implementation in torch. [Ex : VGGNET, ALEXNET]

Torch was good but it had mainly 2 problems :-

1) Torch was LUA [ programming Lang ] based framework.

So if we wanted to build any app using torch then we had to write code in LUA.

[ since LUA was not a popular lang , therefore a beginner has to first learn Lua, then torch , so there was a bit friction ]

So this friction was recognised by Meta AI researchers , they knew torch was very powerful & optimised for tensor based operations. But we had to code in LUA .

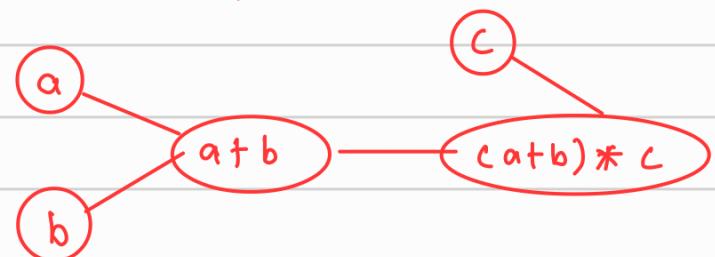
etnd the whole ML/DL community codes in Python. Therefore they thought what if we merge these two ideas. Meaning create a new library where all the powerful capabilities of torch. was present [ tensor calculation ] & we code it in python. So from this idea came ↴

Pytorch.

torch + Python = Pytorch.

# Pytorch 0.1 2017 [ Version 1 ]

- first version.
  - [ 2 core Features ]
- 1) Python compatibility → Even compatible with other libs of python like numpy , scipy etc .
- 2) Dynamic Computation Graph
  - ↳ visual way of representing mathematical operations.



neural networks are also after all mathematical functions .

Ex :  $\sigma(Wx + b)$  → graphs → computation graphs .

deep learning libraries like pytorch & tensorflow etc internally store these operations using graph & these graphs are called computation graphs .

2nd problem of torch → computation graphs were static

once we build our model , a comp graph is generated & that comp graph is fixed can't make any changes . training will now be done on that only .

Static Computation graph once made → can't be changed.

- ↳ flexibility reduces.
- ↳ debugging difficult
- ↳ Experimentation is difficult.

Pytorch identified the problem & rather than going for static they made dynamic Comp graph.

↳ can be changed at runtime (not fixated)

- ↳ flexibility in system.
- ↳ debugging easier
- ↳ Experimentation easier.

Ex Analogy :

Gov agency → making roads.

↓  
before making roads

Scenario 1 :

we make a blue print.



blueprint final → road will be constructed based on it  
↓

Example of static Comp graph.

Scenario 2 :



when road was about to be constructed that there is temple in b/w & blueprint shows road will be built on it. obviously we can't demolish temple & build road across it.  
∴ we need to change.



Ex of dynamic comp graph.



. can change on the go during runtime.

Researchers loved this lib [since Experimentation was easier.] ∴ there was a widespread adoption among researchers



adv: future research papers mentioned pytorch a lot.  
(pytward) [Code was written in pytorch]

## Pytorch 1.0 [ 2018 ]

Although pytorch was popular among research , adoption in industry was not that much [ in production & deployment tensorflow was used more ]

↳ noticed by pytorch team . → in next release ( 2018 ) they focused on this .

2 major changes :

→ torchscript was introduced .

↳ used for model serialisation & optimisation .

when we develop a model using Pytorch then to run it we use pytorch as well but if we serialise our model , then we can send our serialised model to somewhere where there is no pytorch or even python , then also we can run the model .

Ex : developing an android app → want to add dl feature  
[ no <sup>↑</sup> python & pytorch ]

If we use pytorch during training & using torchscript we serialise it , then the serialised model can be run on the android app .

at this time meta researchers were working on another dl lib called Caffe 2.

{

power dl lib

specially used for big models  
& big data set.

[ high performance scalable dl library ]

[ used in production ]

[ production grade deployment capabilities ].

{ Pytorch → good for research

Caffe 2 → good for production grade deployment

They merged Caffe 2 into Pytorch.

after this merger Pytorch became good in production as well.

## Pytorch 1.x Series



Support for distributed training.

[multiple machines, multiple GPUs.]



ONNX [open neural network Exchange] support.



open source format.

till now, pytorch made model can be used in Pytorch  
tensorflow made model can be used in tensorflow.

Basically interoperability b/w lib's was missing.

with the introduction of ONNX, if we save our model  
in ONNX format then our trained model can be  
used in other lib's as well. [for prediction, inference etc]



Quantisation introduced.



our dl [models] is → [weights] are basically → [nums]

↓  
and numbers are  
stored in a data

type for ex:

takes [32 bit] ← [float]  
memory.

Let's say our model have 1 million weights so  
space occupied by our model  $\underbrace{32 \text{ million bits}}_{\text{model size}}$ .

So let's say we want to deploy it in an android app.  
we don't want to unnecessarily fit this many size onto  
our app so ∵ we need to apply some techniques to reduce  
the size of our model.

One of those technique is Quantisation

→ rather than storing our  
weights in bits say float  
we use a lower precision data  
type bits say integer.  
now integer uses 8 bits.

∴ 8 bits  $\times$  1 million.

8 million bits.

→ size of model reduced.

This is model Quantisation. → used for compression  
of model.

- ↳ Ecosystem libs were introduced → like for vision domain → torch vision (CV)
  - ↳ we will find a lot of datasets related to CV
  - ↳ a lot of utilities func
  - ↳ preprocessing func.
  - ↳ pretrained model.

Similarly for nlp domain → torchtext (NLP)

also for audio → torchaudio (audio)

Impact: as pytorch improved → community expanded.

↓  
contribution increased  
↓

community based libs introduced.

high level api      Ex ⇒ Pytorch lightning  
built on top of      ↳ keras for pytorch.  
pytorch.              ↳ api for tf

⇒ hugging face transformers  
[ initial versions were  
written in pytorch ] .

→ as dl community started using pytorch,  
cloud services introduced native pytorch support .

## Pytorch 2.0 [ Version 2 ]

⇒ optimisations .

    ↳ latency → less time more predictions / performance .

    ↳ throughput → in one time how many data we can process

⇒ Improved compilation techniques .

⇒ Optimised the running of pytorch on hardwares . [ TPUs ,  
custom AI chips  
etc .... ]