



IT UNIVERSITY OF COPENHAGEN

www.dasya.dk @dasyaITU https://daphne-eu.github.io/

https://rad.itu.dk

<u>www.itu.dk</u>

peaceful co-habitation on GPUs for deep learning

Pınar Tözün

Associate Professor, IT University of Copenhagen

pito@itu.dk, www.pinartozun.com, @pinartozun

Microsoft GSL Talk Series

24/05/2022

unsustainable growth of deep learning

2012 present

- powerful hardware
- larger datasets
- deep learning frameworks



300000x increase in computational need

for deep learning models.

- computational efficiency is ignored
- → main performance metric = *accuracy*
- high computation (carbon) footprint
- \rightarrow ... with low transparency
- throw new & expensive hardware at the problem?
- ➔ no, there is no free lunch

how did I get into deep learning?

sebastian baunsgaard

IT UNIVERSITY OF COPENHAGEN

year 2018

me



Could you supervise our MSc thesis?

What would you like to work on?

Automatic speech recognition

Why are you talking to me?

We want to make it scalable

ok then



how did I get into deep learning?

sebastian baunsgaard

IT UNIVERSITY OF COPENHAGEN

year 2018

me



Could you supervise our MSc thesis?

What would you like to work on?

Automatic speech recognition

Why are you talking to me?

We want to make it scalable

ok then



sebastian benjamin wrede



agenda

- training speech recognition on co-processors
- studying workload co-location

• challenges & opportunities

speech recognition



- human-computer & human-human interactions
- hospitals, call-centers, virtual assistants, etc.

state-of-the-art *acoustic models* are based on neural networks in recent years → natural fit for GPUs

acoustic model

input features
$$\longrightarrow$$
 CNN $\begin{bmatrix} 5x \\ LSTM \end{bmatrix}$ FFNN \longrightarrow character probabilities

- inspired by Baidu Research, Deep Speech 2, ICML 2016
- basis for MLPerf's speech recognition benchmark as well

input features
$$\rightarrow$$
 $1x-3x$ $1x-7x$ $FFNN \rightarrow$ character probabilities deep speech 2

process of determining the right set of layers is heavily based on trial-&-error

sebastians built *rebelrig*



motherboard = repurposed cheap crypto-mining rig





experimental setup

• acoustic model implemented using TensorFlow 1.14

 training over three hardware platforms

- dataset : LibriSpeech
 - audiobooks
 - ~1000 hours of speech (both clean & noisy)



price/performance results



no huge difference in accuracy across platforms

high throughput != faster time-to-accuracy

impact of batch size



lessons learned

- very powerful co-processors more and more widely available for machine learning / deep learning
- but takes a lot to exploit, no free lunch as usual
- need to invest further in improving machine learning libraries and hardware resource managers
- on the other hand, low-budget platforms may be good enough for your needs

same old challenge for data-intensive systems, different workload & hardware

agenda

- training speech recognition on co-processors
- studying workload co-location

• challenges & opportunities

workload co-location

multiple workloads sharing hardware resources

benefits when a single workload cannot utilize available resources

main high-level challenge = interference across workloads

workload co-location on (NVIDIA) GPUs

• vanilla co-location

kernels of different workloads are time-multiplexed (not concurrent)

• virtualization

• practical, but also based on time-sharing

• multi-process service (MPS)

- GPU resources are split (auto-magically) across co-located workloads
- kernels of different applications can run simultaneously
- allowed for one user (for safety reasons)

• multi-instance GPU (MIG)

- hardware support for resource split, introduced with NVIDIA A100
- can do all of the above in a MIG partition





unavailable compute unit

Х







Х



unavailable compute unit = less than 14 SMs

performance impact of MIG-based co-location

NVIDIA DGX Station A100



CPU = AMD 7742 – 512 GB RAM 64 physical cores GPU = NVIDIA A100 – 40 GB RAM allows *multi-instance GPU (MIG)*

TensorFlow^{2.7}

workloads	model	dataset
small	ResNet26	CIFAR-10
medium	ResNet50	downsampled ImageNet <u>*</u>
large	ResNet152	ImageNet (2012)

batch size = 32 for all runs on single GPU

- 25 epochs for small
- 5 epochs for medium & large

MSc thesis work of Stilyan Petrov Paleykov & Anders Friis Kaas



time per epoch – small case



opportunity to co-locate training runs with slight latency increase

time per epoch – medium case



configurations

1g.5gb case isn't feasible anymore due to insufficient memory, not much gain from co-location

time per epoch – large case



similar to medium case in terms of co-location overall, parallel runs don't interfere as long as there is enough memory

IT UNIVERSITY OF COPENHAGEN

GPU utilization



fine-grained parallel runs increase utilization for *small case medium & large* cases utilize the whole GPU well without parallel runs

challenges & opportunities

opportunities

- GPUs that allow finer-grained scheduling & space management
- diversity of applications, hardware, & end-users
- ➔ creates opportunity for effective resource sharing on GPUs

challenges

- representative workloads
- experimental duration
- profiling & co-location granularity

team **RAD** – resource-aware data systems

IT UNIVERSITY OF COPENHAGEN





Ties Robroek

Ehsan Yousefzadeh-Asl-Miandoab



Jon Voigt Tøttrup

edge / IoT devices



Robert Bayer



Lottie Greenwood

IT system admin

servers with CPU-GPU co-processors

https://rad.itu.dk

challenges & opportunities

thank you!

opportunities

- GPUs that allow finer-grained scheduling & space management
- diversity of applications, hardware, & end-users
- Creates opportunity for effective resource sharing on GPUs

challenges

- representative workloads
- experimental duration
- profiling & co-location granularity