

# GOTO AARHUS 2022





#### WHO AM I



#### Zenodia Charpy

#### Senior Deep Learning Data Scientist @ NVIDIA - NeMo| Megatron | bigNLP

My past experience:

Azure : senior data scientist & solution architect

+8 years experience as in-house data scientist & data science consultant

# 2GPT3, NLP AT BIG SCALE

### Overview

The 2 GPT3s

The big 530B Megatron-Turing NLG GPT3 model

Motivation - Why Very Large Language Models

3 basic ingredients

Steps to follow - workflow

QnA – reach out and kick-start

### The 2 GPT3

One *English*, the other *Swedish* 

	English	Swedish
Data size	Terabytes of data	~40GB
Model size	530B	760M
Compute	<b>4480</b> GPUs (560 nodes)	<b>8</b> GPUs (1 nodes)
Time to train	~ <b>20</b> days	~7 days
Use cases	Mutiple downstream tasks with 0/1/fewshots	Limited to a handfull of cases

		English	Swedish
9	Data size	Terabytes of data	~40GB
0	Model size	530B	760M
	Compute	4480 GPUs (560 nodes)	16 GPUs (2 nodes)
	Time to train	~20 days	~7 days
	Use cases	Mutiple downstream tasks with 0/1/fewshots	Limited to a handfull of cases
530B	MEGATRON	TURING NLG	

- ENGLISH GPT3

# VIRTUAL ASSISTANCE





#### Customer Service Kiosk with NVIDIA Omniverse Avatar for Project Tokkio

#### Large NLP models powers:

 $\bigcirc$  Zero-shot intent and slot classification

- Context-based extractive Q&A
- Conversational phrasing of answer

# TOY JENSEN - ASK ME ANYTHING





#### Expert, Natural Q&A

with NVIDIA Omniverse Avatar for Project Tokkio



# OTHER USE CASES

#### SOLVE RIDDLE

 $Context \rightarrow$  Here is a riddle:

All of us have one, but few get to choose If you don't know mine, you are not my friend When it is called, our attention is drawn Even if we are gone, they are still around

I think the answer is

#### **CODE GENERATION**

 $\texttt{Context} \rightarrow$ 

```
def find_3_or_7_divisible_fibs(n):
    """Find all the Fibonacci numbers below n that are divisible by 3 or divisible by
    7.
    """
```

WHY VERY LARGE LANGUAGE MODELS

### WHY **PRE-TRAIN** VERY LARGE NLP MODELS

Train Once, Perform different tasks

Step 1: Train a *Very Deep/Huge* model

General

Knowledge

Science



Step 2. Query different tasks with *a few labelled* data

#### Huge means Billions of parameters

Literature

#### MODEL SIZE GROWS EXPONENTIALLY 1000 Model Size (in billions of parameters) GPT-3 **Megatron-Turing** (175B) NLG (530B) 100 Megatron-LM Turing-NLG (8.3B) (17.2B) 10 T5 (11B) GPT-2 1 (1.5B) BERT-Large (340M) 0.1 ELMo (94M) 0.01 2019 2020 2021 2018 2022

Link to the article





BUT THAT'S ENGLISH
- WHAT ABOUT LOCALIZED LANGUAGE ?

# **DVIDIA**.

	English	Swedish
Data size	Terabytes of data	~40GB
Model size	530B	760M
Compute	4480 GPUs (560 nodes)	16 GPUs (2 nodes)
me to train	~20 days	~7 days
Jse cases	Mutiple downstream tasks with 0/1/fewshots	Limited to a handfull of cases

SECURE 3 BASIC INGREDIENTS

#### **3 BASIC INGREDIENTS**

#### **The Data**

#### The Compute

#### **The Framework**







WHO HAS ALREADY BENEFITTED FROM THIS WORKFLOW ?

### WHO HAS ALREADY USED THIS WORKFLOW ?

And trained + deployed their GPT3 models





Jun 8 · 3 min read · 🕑 Listen

As the name suggest, **GPT-SW3** is the Swedish counterpart to GPT-3: a largescale generative pretrained Transformer built on massive amounts of Swedish data with the explicit purpose of being able to generate Swedish text, and to thereby function in zero-shot and few-shot scenarios. The initiative to develop such a model for Swedish is led by the <u>NLU research group</u> at <u>AI Sweden</u>, and includes a number of collaborators, including <u>RISE</u>, the <u>WASP WARA on media</u> <u>and language</u>, and <u>Nvidia</u>. GPT-SW3 already exists in a preliminary version with 3.5 billion parameters, which was trained on a moderately sized corpus of approximately 100 GB of mostly web data. Our current goal is to train a 175 billion parameter model with close to 1 TB of data.

#### GPT-3 på svenska – de bygger en "superintelligens" som pratar svenska

2022-06-09 06:00 Av: Simon Campanello 0 kommentarer

Link to the article



Applicable to any language



Scalable

KICK START YOUR JOURNEY - THE WORKFLOWS

#### THE WORKFLOW - START THE JOURNEY

Generic workflow - works for small to large to very large NLP models



1.Problem sizing



### THE ITERATION TIME

Problem sizing - reducing years to weeks/days

GPT-3		Estimated Tin	ne to Train		
Model Parameter Count	5 x SuperPod 100 nodes (days)	3 x SuperPod 60 nodes (days)	1 SuperPod 20 nodes (days)(weeks)	4 x DGX A100 (weeks)( <mark>years</mark> )	
0.5B	0.21	0.21	0.62	0.44	
18B	5	8	23	16	
175B	43	72	31	3.0	
		Training			

Link to python function to estimate training time based on the equation from the paper below

Deepak Narayanan et al. Efficient Large-Scale Language Model Training on GPU Clusters Using Megatron-LM.

NVIDIA DGX A100 SuperPOD



2.Collect data - part 1/3- Training data



### DATA COLLECTION

#### Training data: Quantity vs. Quality







Figure 1: Treemap of Pile components by effective size.



An 800GB Dataset of Diverse Text for Language Modeling <u>https://pile.eleuther.ai/</u>

### WEIGHING AGAINST QUALITY VS QUANTITY

2.Data Collection

Weighing mechanism is specific to the datasets used for training

Dataset	Tokens (billions)	Weights (%)	Epochs
Books3	25.7	14.3	1.5
OpenWebText2	14.8	19.3	3.6
Stack Exchange	11.6	5.7	1.4
PubMed Abstracts	4.4	2.9	1.8
Wikipedia	4.2	4.8	3.2
Gutenberg (PG-19)	2.7	0.9	0.9
BookCorpus2	1.5	1.0	1.8
NIH ExPorter	0.3	0.2	1.8
Pile-CC	49.8	9.4	0.5
ArXiv	20.8	1.4	0.2
GitHub	24.3	1.6	0.2
CC-2020-50	68.7	13.0	0.5
CC-2021-04	82.6	15.7	0.5
RealNews	21.9	9.0	1.1
CC-Stories	5.3	0.9	0.5

2.Collect data - part 2/3- Data for downstream tasks



### DATA COLLECTION

#### Annotating and creating data for downstream tasks

Train the GPT model unsupervised

Evaluate





### DATA + <u>LABEL</u> COLLECTION

#### Annotation requires Tools and Expertise

UNIVERSITY OF GOTHENBURG	svenska RÅKBANKENTEXT
NEWS AND EVENTS BLOG	RESEARCH TOOLS RESOURCES FAQ ABOUT US CONTACT US
<u>Home / Tools</u> / <u>Sparv</u> / An	notations by Sparv
Sparv	Annotations by Sparv
Sparv Pipeline Sparv's user manual	This page aims at giving an overview of the analyses available within the Sparv Pipeline and the Sparv plugins developed by Språkbanken.
Annotations by Sparv	annotations are the names of the annotations as they are used in the corpus config file in the export.annotations section (read more about this in the <u>Sparv user manual</u> ). Please observe that the annotations usually have shorter names in the corpus exports
Web Sparv	annotators are the names of the annotation functions (including their module names) which are used for procuding the annotations. They can be run directly with the sparv run-rule [annoterare] command. Though in most cases this is not necessary due to the fact that the annotation functions needed to produce all annotations listed in the corpus config file are executed automatically when running sparv run.
	Analyses for contemporary Swedish
	The following analyses are available in Sparv for contemporary Swedish. For various reasons not all of these are used by us at Språkbanken for analysing the corpora in Korp.
	Sentence segmentation with PunktSentenceTokenizer
	<ul> <li>description: Texts are split into sentences.</li> <li>model: <u>punkt-nltk-svenska.pickle</u></li> <li>method: The model is build with <u>NLTK's PunktTrainer</u> and trained on <u>StorSUC</u>. The segmentation is done with NLTK's <u>PunktSentenceTokenizer</u>.</li> <li>annotations:         <ul> <li>segment.sentence: sentence segments</li> </ul> </li> </ul>
	annotators: segment.sentence

#### Stockholm-Umeå corpus 3.0

The Stockholm-Umeå Corpus (SUC) is a collection of Swedish texts from the 1990's, consisting of one million words in total. The corpus is balanced, meaning that it contains various text types and stylistic levels. The texts are annotated with part-of-speech tags, morphological analysis and lemma (all that can be considered *gold standard data*), as well as some structural and functional information.

*Version 1.0* was developed in co-operation between Gunnel Källgren at Stockholm University and Eva Ejerhed at Umeå University and was made available in 1997 by the department of linguistics at Stockholm University. *Version 2.0* was made available in 2006 by Sofia Gustafson-Capková and Britt Hartmann at the department of linguistics at Stockholm University. It contains the same texts as SUC 1.0 but is extended with some annotation. Additionally, SUC 2.0 contains bonus materials. TigerSUC is SUC 2.0 converted to TIGER-XML by Martin Volk. StorSUC is additional SUC material of four million words. *Version 3.0* is available since 2012. It contains improved annotations, and unannotated texts with seven million words. (For the TigerXML-version, Suc2c, Suc2d, and the DTDs we still refer to version 2.0.)

2. Collect data - part 3/3- Legal/ethical/responsibleAl concerns



#### DATA COLLECTION Establish Censoring Criteria



- 1. Data source licence implications
- 2. Data characteristics and coverage
- 3. Sufficiency of data pre-processing pipeline
- 4. Addressing issues of dataset bias
- 5. Addressing issues of dataset fairness
- 6. Addressing to ethical concerns
- 7. Addressing to legal concernss
- 8. Blacklisting & toxicity removal
- 9. Composition and weighing
- 10. ...

3.Data preprocessing - part 1/3- Filter / clean / deduplicate



#### DATA PRE-PROCESSING

#### Filtering + Deduplicate + Sentence Boundary

s_ roo	t@a0ca	a60d991fa: /works1X Ada_cleaning_demo.ipynb	
8	+ %	The the second s	thon 3 C
		Step 1 - detect and filter undesired langauge in the raw text corpus	
	[1]:	<pre>import random from langdetect import detect f=open('./bland_langs.txt','r') lines=f.readlines() len(lines)</pre>	
	[1]:	1824	
	[]:	<pre>rn=random.randint(0,len(lines)-1) print("sentence = ", lines[rn]) lang=detect(lines[rn])  print("detected language is = ", lang)</pre>	
		step 2 - deduplicate based on similarity threshold	
		sanity check, fetch previously labelled two documents with groundtruth (i.e duplicated or not)	
		check against threhold, set to 0.85 and determine these two documents are duplicated or not !	
	[]:	<pre>import itertools from lsh import cache, minhash # https://github.com/mattilyra/lsh</pre>	
		<pre># a pure python shingling function that will be used in comparing # LSH to true Jaccard similarities def shingles(text, char_ngram=5):     return set(text[head:head + char_ngram] for head in range(0, len(text) - char_ngram))</pre>	

Demo: Data pre-processing toy example for the Swedish language

Languages filtering + Deduplication + Sentence boundary (filter short sentences/single sentence document)

Link to data cleaning example notebook

3.Data preprocessing - part 2/3- Train your own tokenizer

### TRAIN OWN GPT TOKENIZER

3.Data Pre-processing

#### Train a GPT tokenizer on your own raw corpus

s root@a0ca60d991fa: /works ×	≣ trainGPTTokenizer.py ×
root@a0ca60d991fa:/workspace/SVdata/gtc# python trainGPTTokenizer.pyinfil	<pre>parser.ado_argument(pretrained_gpt_oir , default=None, type=str, netp= path to pretrained or vocab and merge files, default None') parser.add_argument(' vocab size' default=None type=int</pre>
	19 help='specify the vocab size when training HF GPTBPE for own language usually 16k/32k
	/48k/64k')
	20 args = parser.parse_args()
	<pre>21 tokenizer = Tokenizer(models.BPE()) 22 if and protocols.BPE())</pre>
	17 args.toad_pretrained :
	$23$ print("loading gpt_bpe english vocab and merge \n")
	25 vocab file=args.pretrained opt dir+'/opt2-vocab ison'
	26 merge file=args.pretrained gpt dir+'/gpt2-merges.txt'
R	<pre>27 tokenizer.model = models.BPE.from_file(vocab_file, merge_file)</pre>
L. L	28 else:
	<pre>29 print("please provide path to the pretrained gpt vocab and merge file !")</pre>
	30
	<pre>31 tokenizer.pre_tokenizer = pre_tokenizers.ByteLevel() 23 tokenizer deceder = PutelouglDeceder()</pre>
	33
	34 print("include minimal special token end of text ")
	35 special tokens= [" <lendoftext]>"]</lendoftext]>
	36
	37 <i># Set the training hyperparameters</i>
	<pre>38 trainer = trainers.BpeTrainer(</pre>
	<pre>39 vocab_size=args.vocab_size,</pre>
	40 special_tokens=special_tokens,
	41 Initiat_atphabet=pre_tokenizers.ByteLevet.atphabet()
	43 # Train it with either files or an iterator:
	44 tokenizer.train([args.infile], trainer=trainer)
	<pre>45 print("Trained vocab size: {}".format(tokenizer.get vocab size()))</pre>
	46 <i># You will see the generated files in the output.</i>
	<pre>47 print("saving trained BPE model to : ", args.bpe_path)</pre>
	48 tokenizer.model.save(args.bpe_path)
	49 print("model saved ! \n\n\n")
	50 print("testing\n\n\n") 51 test tyt="Her percent for for the state of the second
	som frågat Thermia? Skulle det inte vara väsentligt att kunna kolla historiken nå den då man skall ställa in

Demo: Train a GPTBPE Tokenizer on a toy Swedish data Link to <u>train GPTBPE Tokenizer example notebook</u>

3.Data preprocessing - 3/3- Optimal data loading



### PREPROCESS TO BINARY FILES

Improve data loading performance during training time

<u>n inter anna ana anna anna anna cointe</u>	1/26.0
+ 🗈 ±	C
/ dataset /	
ame	Last Modifie
ld_ckpt	11 days ag
∎ zh_glue	11 days ag
bk_text_sentence_test_indexmap_80mns_512msl_0.10ssp_1234s.npy	23 days ag
bk_text_sentence_train_indexmap_8000000mns_512msl_0.10ssp_1234s.npy	23 days ag
bk_text_sentence_valid_indexmap_800080mns_512msl_0.10ssp_1234s.npy	23 days ag
Bk_text_sentence.bin	23 days ag
B bk_text_sentence.idx	23 days ag
<u>a demo_gpt_sentence_text_sentence.bin</u>	an hour ag
mmap:	
process data into xxx.idx + xxx.bir optimized memory mapping to imp performance	n for prove

4. Training considerations



### CONSIDERATION FOR TRAINING LAUNCHES

Manual vs automatic : obtain optimal training/inference configurations

DIY(Megatron-LM)



- 1. Manual runs to obtain optimal training config
- 2. Manually develop inference route for pretrained models
- 3. <u>Manually</u> debug with limited support
- Control



Manual work



### Enterprise ready(NeMo Megatron)



- 1. Get <u>automatic</u> recommandation of optimal training config from the HP tool
- 2. Support by default training and Inference per model size in the template
- 3. Get enteprise-grade expertise support

ଡ଼ଡ଼ଡ଼



### **DIY**- MANUALLY TUNE TO ENSURE PERF VIA PROFILING

Profile training runs | improve iteratively



4.Language Model Pretraining

### **DIY** - OBTAIN GOOD CONFIG FOR MODEL PRETRAINING

4.Language Model Pretraining

#### Training configuration performance obtained via manual experiments

based on 100		Activation	4						312	A100 FP16 TFLC	PS		
		No Activation	3						112 V100 PCIe FP16 TFLOF				
elapsed time per					ind goo	d trainin	g config		125	V100 SXM			
Speed ms	N Gpus	Number of Nodes	Model	Micro	DP Size	TP Size	PP Size	PP	NVIDIA	Lower bound ut	ilisation		
31579	8	2	18	4	. 1	4	4	128	156.514		50.16%		
44438	8	2	18	6	1	4	4	128	166.835	Teraflop	53.47%		
31532	8	2	18	4	- 1	4	8	128	156,747	util %	50.24%		
44688	8	2	18	6	1	4	8	128	165.902		53.17%		
62078	8	2	18	4	. 1	4	8	256	<del>159</del> ,236		51.04%		
15889	8	2	18	8	1	4	8	32	155.533		49.85%		
16355	8	2	18	4	- 1	4	8	64	151.102		48.43%		
31191	. 8	4	39	4	1	4	8	128	167.118		53.56%		
45146	8	4	39	6	1	4	8	128	173.191		55.51%		
87563	8	4	39	6	1	4	8	256	178.588		57.24%		
173505	8	4	39	6	1	4	8	512	180.257		57.77%		
43119	8	8	52	4	- 1	4	16	256	160.871		51.56%		
43143	8	8	52	4	- 1	4	16	256	160.782		51.53%		
43120	8	8	52	4	. 1	4	16	256	160.868		51.56%		
43154	. 8	8	52	4	. 1	4	16	256	160.741		51.52%		
62427	8	8	52	6	1	4	16	256	166.673		53.42%		
80920	8	8	52	8	1	4	16	256	171.444		54.95%		
121531	. 8	8	52	6	1	4	16	512	171.231		54.88%		
158302	. 8	8	52	8	1	4	16	512	175.275		56.18%		
43112	8	8	52	4	. 1	4	32	256	160.897		51.57%		

### **ENTERPRISE READY** - PRETRAINING WITH RECOMMENDED CONFIG

4.Language Model Pretraining

NVIDIA NeMo Megatron makes training easy with recommended config & perf guarantee



### **ENTERPRISE READY** - PRETRAINING WITH RECOMMENDED CONFIG

4.Language Model Pretraining

NVIDIA NeMo Megatron makes training easy with ready-made template config & perf guarantee



Training Config



### MONITOR TENSORBOARD

ensorBoard SCALARS											INAC	TIVE	- C 🌣	0
Show data download links	<b>Q</b> Filter tags (regu	lar expressions supported	)											
] Ignore outliers in chart scaling	grad-norm	grad-norm												1
poltip sorting method: default 🔹	grad-norm vs sam	ples												1
noothing	learning-rate													1
0	learning-rate vs sa	amples												1
	Im loss													1
contal Axis	Im loss													
RELATIVE WALL														
	12													
a regex to filter runs	10													
Ο,	8													
TOGGLE ALL RUNS	6													
		all and a												
	4		and a low second second	and the second second second			and the second	and and be defined					-	
	2													
	0								1.					_
														_
		0 50k	100k	150k	200k	250k 300	k 3!	50k 40	10k 45	50k	500k	550k	600k	
	Im loss validation													1
	Im loss validation	vs samples												1
	Im loss vs sample	s												1

5. Deploy & Serve- Text generation demo

1.P S	roblem 2.Data izing Collection Pre	3.Data Processing Pretraining 5.Deployment 6.Downstream Tasks
$\leftarrow$	$\rightarrow$ C	○     □     localhost:2345/lab?
0	File Edit View Run H	Kernel Tabs Settings Help
	+ 🗈 ± C	Is root@3c633147b13b: /works × Step3_ViewGenerated_text.i● E Step2_760MGPTGenerateTc ×
	🖿 / SVdata / gpt2bpe /	use contiguous buffers in ddp False
0	Name 🔺	use_cpu_initialization None
	🖿 32k	use_one_sent_docs False
-	🖿 48k	virtual_pipeline_model_parallel_size None
6	56k	vocab_file
	□ apt2-merges txt	weight_decay 0.01
<sup>b</sup> o	(1) ant2-vocab ison	world_sizeend of arguments
	SVGPT 32k text d	setting number of micro-batches to constant 1
-		> building GPT2BPETokenizer tokenizer
_	SVGP1_32K_lext_d	> padded vocab (size: 32000) with 0 dummy tokens (new size: 32000)
		> initializing tensor model parallel with size 1
*		> initializing pipeline model parallel with size 1
		> setting random seeds to 1
		> initializing model parallel cuda seeds on global rank 0, model parallel rank 0, and data parallel rank 0 with model parallel seed: 2719 and data parallel seed: 1
		make: Entering directory '/workspace/megatron/data'
		make: Nothing to by done for 'default'.
		make: Leaving directory '/workspace/megatron/data'
		>>> done with dataset index builder. Compilation time: 0.038 seconds
		> compiling and loading fused kernels
		Detected CUDA files, patching ldflags
		Emitting ninja build file /workspace/megatron/fused_kernels/build/build.ninja
		Building extension module fused mix prec layer norm cuda
		ninia: no work to do.
		Loading extension module fused_mix_prec_layer_norm_cuda
		>>> done with compiling and loading fused kernels. Compilation time: 1.028 seconds
		Dullding GPT model
		loading checkpoint from /workspace/tools/gpt3 iter2mil/ at iteration 2000000
		checkpoint version 3.0
		successfully loaded checkpoint from /workspace/tools/gpt3_iter2mil/ at iteration 2000000
		Avg s/balch: 0.90//0340/0098140

5.Deploy & Serve



#### Enterprise ready(NeMo Megatron) -

support by default *automtaic* deployment route





1.Problem<br/>Sizing2.Data<br/>Collection3.Data<br/>Pre-processing4.Language Model<br/>Pretraining5.Deployment6.Downstream<br/>Tasks

P-tuning for downstream tasks

### NLP APPROACH (CIRCA 2019)

Fine tune per specific task

Step 1: Pre-training a Encoder

Step 2. Fine tune for a specific task







### NEW NLP APPROACH (CIRCA 2021)

#### Train Once, Perform different tasks

Step 1: Train a *Very Deep/Huge* model

General

Knowledge

Science



Step 2. Query different tasks with *a few labelled* data

#### Huge means Billions of parameters

Literature

GitHub

#### P-TUNING METHOD Vanilla P-tuning



### **BRING YOUR OWN DATA VIA P-TUNING**

Well-engineered architecture in NeMo (future release in NeMo Megatron)





### BUT HOW DO I TRY-OUT/KICK-START FOR REAL?

# MULTIPLE WAYS TO TRY-OUT/KICK-START

#### Request a bootcamp

← → C ( a github.com/openhackatho	ons-org/gpubootcamp/tree/master/ai/Megatron	ピ ☆ 🗯 🖬 📀
pytorch 🚺 Thomas_Mixed_pre 🚺 Nicola	1_share 💧 Nvidia- Google Drive 🤕 transferlearningToo 🎯 NvidiaTransfi	erLearn 🧕 Scenarios and Acci 🎧 models/model_zoo » 📔 Other book
G openhackathons-org / gpub	pootcamp Public 🛇 Ed	fit Pins ▼ ③ Unwatch 17 ▼ ¥ Fork 159 ☆ Star 266 ▼
↔ Code ⊙ Issues 28 \$1 Pull re	equests 3 🖓 Discussions 🕑 Actions 🖽 Projects 3 🛛	뙤 Wiki ① Security 🗠 Insights
1 <sup>3</sup> master • gpubootcamp / ai /	Megatron /	Go to file Add file * ····
mozhgan-kch megatron readme up	idated	921966a on Apr 1 🕄 History
English/Python	implement proof-reading thanks to Millie T	8 months ago
C README.md	megatron readme updated	2 months ago
i≣ README.md		Ø

#### Practical Guide to Train Megatron-LM with your own language

Megatron-LM is used in NeMo Megatron, this is the framework to help enterprises overcome the challenges of building and training sophisticated natural language processing models with billions and trillions of parameters.

https://github.com/openhackathonsorg/gpubootcamp/tree/master/ai/Megatron

#### Register LaunchPad



Train a Large-Scale NLP Model with NeMo Megatron

#### Audience: Al practitioner

Products: NVIDIA Base Command, NVIDIA DGX Technologies: NeMo Megatron, PyTorch

GET STARTED >

https://www.nvidia.com/en-us/data-center/launchpad/

#### Join NeMo Megatron EA NeMo Megatron Early Access

NeMo Megatron is a new capability in the NeMo framework that allows developers to effectively train and scale language models to billions of parameters. With NeMo Megatron, you can train different variants of GPT-3 models and scale them to multiple nodes on superpods. This deep learning software stack is optimized for NVIDIA DGX A100 SuperPODs using NVIDIA InfiniBand to provide efficient onpremises compute for training and inferring complex workloads.

Early access to NeMo Megatron is limited to enterprises that want to train and deploy GPT-3 style models on NVIDIA A100 SuperPOD to perform zero-shot tasks such as answering deep domain questions, translating languages, comprehending and summarizing complex documents, etc.

https://developer.nvidia.com/nemo-megatron-early-access

SUMMARY



📀 NVIDIA.

### **REACH OUT**



### Rasmus Bisgaard rbisgaard@nvidia.com



### Zenodia Charpy zcharpy@nvidia.com





# DON'T FORGET TO RATE THE SESSIONS #GOTOaar

Rate a minimum of **5 sessions** and claim your **reward** at the Registration Desk at the Trifork Hall