

NVIDIA Solutions architect

Augmenting LLMs

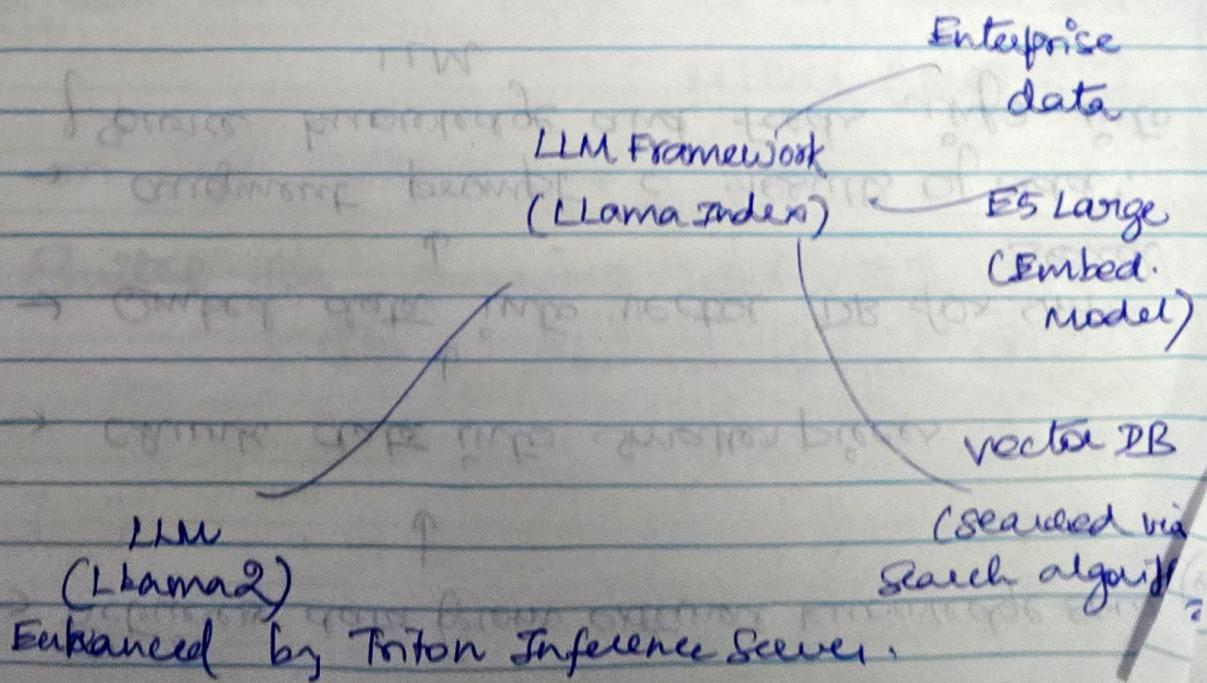
Canonical RAG workflow:

RAG vs Fine-tuning :

- can be combined to leverage strengths
- FT: changes style/tone of response
add new data - impossible / difficult

RAG = Facebook (2020) → generator + I/P + chunk query

GPT : *



NVIDIA Solutions architect

→ Augmenting LLMs using ...

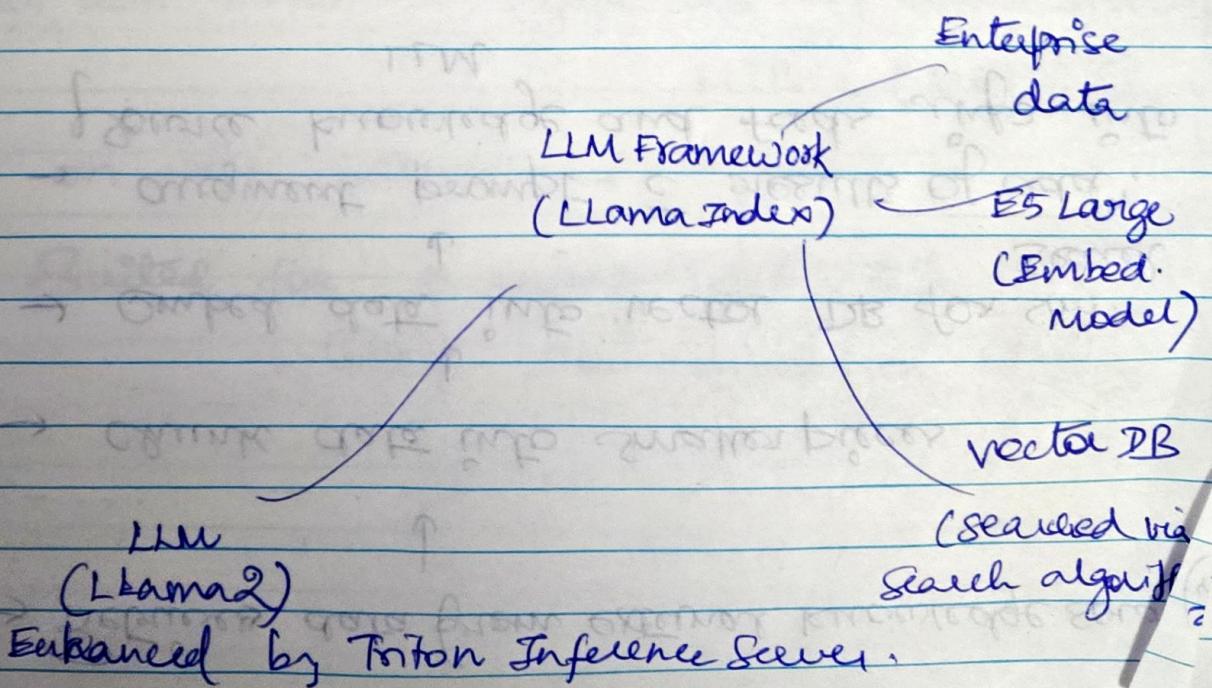
Canonical RAG workflow:

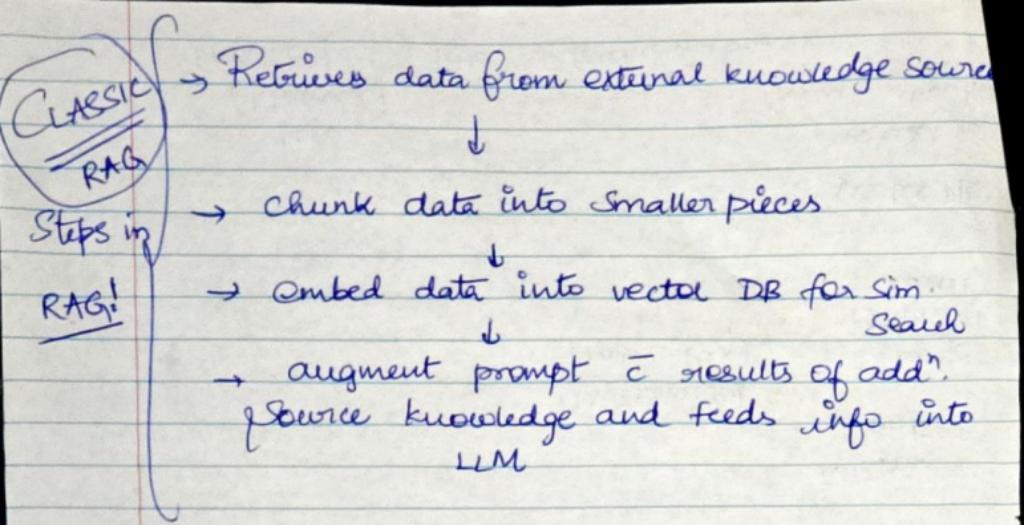
RAG vs Fine-tuning :

- can be combined to leverage strength
- FT: changes style/tone of response
add new data - impossible / difficult

RAG = Facebook (2020) → generator → I/P + chunk query

GPT : ~~✓~~





Agent-assisted RAG :-

COT based RAG : with more toolcalls

Decompose complex query into simple ones

RAG C guardrails :-

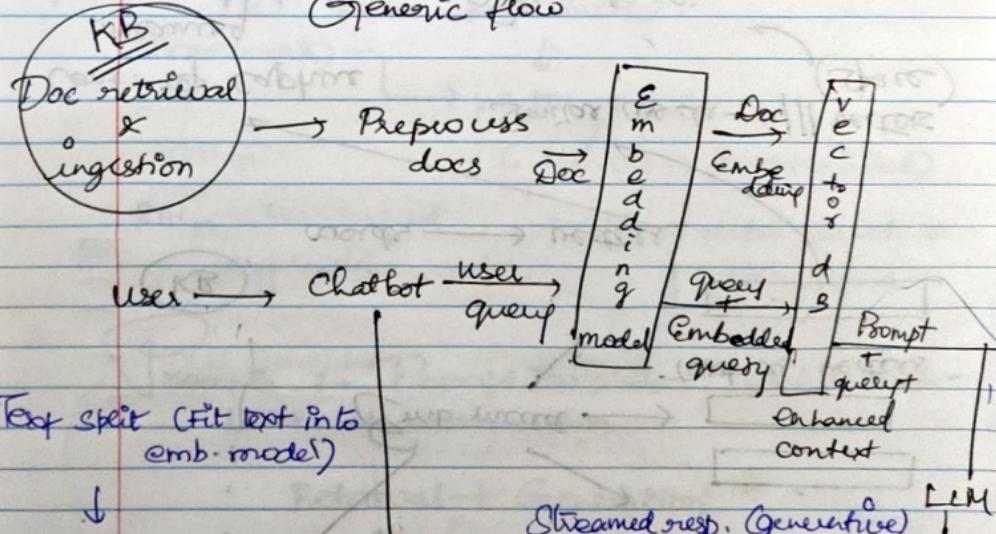
Speeds up decision making process (CoT) when using agents.

→ Outputs align with high quality

guardrails :- a)

b) safety guardrails - filter unwanted content
c) security "

Generic flow



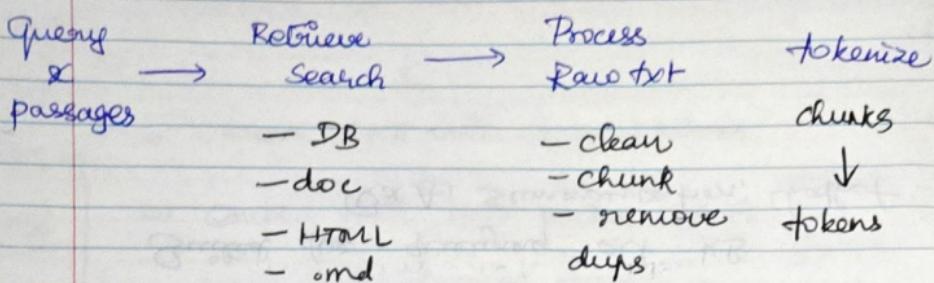
data → High dim vectors → vector DB → rapid search
Milvus (RAPIDS)

Trained on vast datasets →

Suited for querying ext. KB

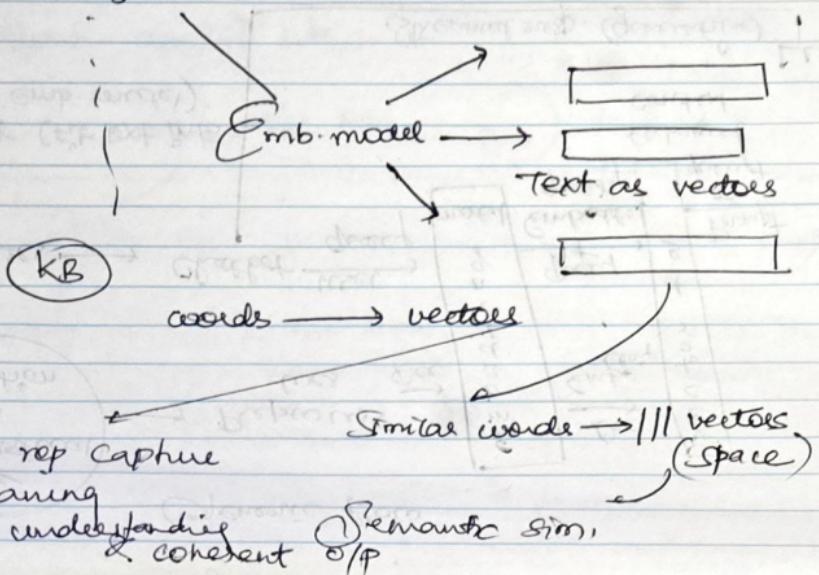
Q&A, Summarization, Chatbot

Ingestion



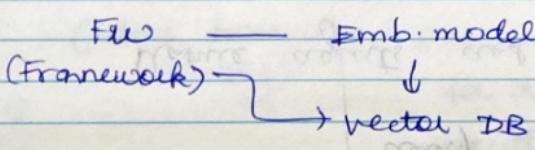
Embedding generation

Ingested data



In-memory / vector store

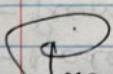
Vector DB → high dim vec.



C libraries → purpose filled

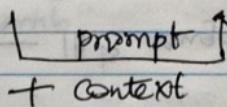
* RAG → Pre-trained Lang. model

Redis → vector search is addⁿ feature.



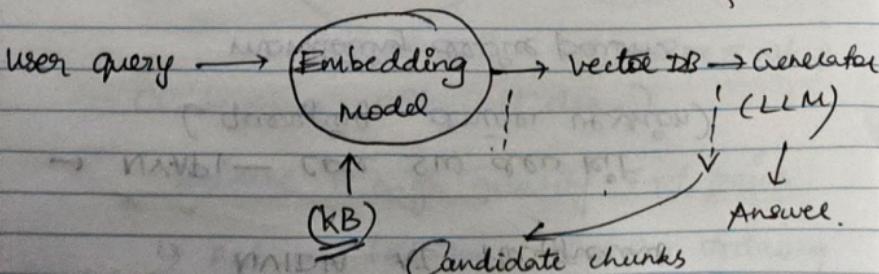
Prompts :- Prompt Template - useful high-level abstraction

Retrieval + Generation



NVDA + RAG : Soln. local inference

Latest Tesla + GTX drivers, Llama 2



NVIDIA AI Playground

→ NVAPI — core SW dev kit
(querying driver version)
monitoring device params

* Interface to Llama 2 API.

Langchain — General framework (matured)

LlamaIndex — API to langchain → Index/query
tasks are good!

Foundation	Chain	Agents
— Model wrapper	combine	— tool-calls
— Prompt template	multiple Comps	to
— Indices / vector store	to solve an	complete task
— memory	issue (chat c)	
— docx Loader/text	pdf, analyze	
— Splitters	data)	

Mature agents and mem. structures

perfor DB → individual DB

LlamaIndex

Data connectors → does data
Indices → Query

(docs from
directory)

Structure data
for diff
usecases

Interface
[Feed in
i/p prompt
& obtain
knowledge
aug. of p]

Reddit/KB

↓

Consistency based
filter (1.3B) Embedding model (E5 Large)

cc-pairs TLM (Transaction Level Modeling)

(2TOM) More eff on GPU — Commercial license!

technique → Contrastive loss

Query emb

Query emb

Pool

Pool

Encoder

Text1

Encoder

Text2

Shared wt.

(SIAMESE N/W)

CC-pairs (Causal clean)

Diverse & high quality txt pairs

↳ Reliable txt pair used for training

Convert
txt to
array of
FP #s

Vector DB: Milvus

* Every vect \longleftrightarrow vect. in collection
DISTANCE

Why Milvus? RAFT accelerator lib. supported

Heterogeneous Comp \rightarrow GPU support,
faster and better perf.

Scaling to more doc \rightarrow better

vector DB integrated C

Chroma, FAISS, MongoDB.

RAPIDS RAFT: Accelerate ML algo on GPU

for
vector search

trained LLM: Meta's Llama?

Stores index & massive vol of emb vectors

[perform better than old models]
SFT: \rightarrow RLHF (Rejection Sampling and PPO)

C opt. Triton (inference server)

ToB params

Why Triton simplifies Inference?

Any framework

GPU+CPU

Any Inf type

(Real type
Batch)
Model Inf.
High Input +
low latency.

CPU/GPU/

Cloud/edge

Designed for

MLOps / webops]

kubernetes, AWS,

GCP AI,

Cons of RAE:

- Tech prof.

Model FT C Nemo.

- Data freshness/retrieval

Guardrails for style

Strategy

- vect search - slow!

RAFT

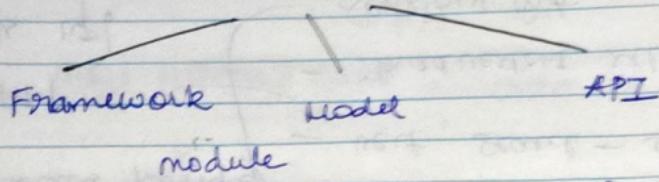
Speeds up!

- Inadvertent release of sensitive info.

Nemo guardrails

Vector DB: Milvus

Nemo : Nvidia Maximum Optimized



- a) Neural model - Speech recognition, language understanding & resp gen
- b) Experiment mgmt - training & eval of Nemo models - perf / bc of Nemo models
- c) checkpointing : Save & load model states during training

Features :-

Mixed precision training (Comb of FP16 & K32 train → ↓ mem usage)

Parallelism > train. → perf model on multi-QPUs

↳ reduces training time & ↑ Scalability

Optimization: Attention optimization, activation recomputation and comm overlap → ↑ model perf and efficiency

use of Nemo :-

chatbot virtual assit speech recog

- VLM →
- a) NeVA (LLava) - text out of img
 - b) LLava-next → PT LLava
 - c) Llama 3.2 vision models : generate text based on imgs

Speech modes @ Diffusion models (High quality speech)

(d) Flux → Speech AI model that generates speech in real-time

Emb models : SBERT, Llama Emb

Tools: Nemo Run → tool to run Nemo models & exp.

PT models - Mixed precision training

* nemo generate --config-type x (Deep learning)

* nemo train --config-type x (Deep learning)

Nemo about Hugging

Nemo agent toolkit

- * nmt run --config-file x --input ...
[Test workflow]
- * nmt serve --config-file x -host ... 127.0.0.1.
[Deploy]

How to write config file?

% cat config.yml

llms:

climate-llm:

-type: nlm

model-name: meta-llm

base-url:

api-key:

temperature:

max-tokens:

workflow:

-type: chat-completion (react-agent)

llm-name: climate-llm

System prompt: |

[no quotes]

Above agent answers generic qns. on climate via LLM data!

Next step:- * Tag Python tools into NAT-compatible tools

* Register multiple tools in React agent workflow

@ register-function(config-type = ...)
def calc-stats(...):

All tools begin in config class

```
class CalcStatsConfig(FunctionBaseConfig,
                      name='calc-stats'):  
    pass
```

```
class CalcStatsInput(BaseModel):  
    country: str = Field(default=' ',  
                          description='...')
```

@register-function(config-type = CalcStatsConfig)

.yaml file :- [project], [project-entry-points.net]

Observability in NAT :-

→ understand what's happening in your workflow

* → Phoenix → Server by Arise

General :

telemetry :

tracing :

phoenix :

-type :

endpoint : (server)

project :

→ Multi-agent integration adding math!

result



→ nat eval -- config-file x



NVIDIA H100: Fast calculator
Data center GPU on Hopper arch. (Better than A100)

Tensor cores: Super-smart math units
to do math problems

- SXM5 PCIe
- ④ recognizes pic and videos
 - ⑤ understand what ppl say
 - ⑥ play games that need math

NVLink - Superfast internet connection
to talk to other comps and share info

A lot of memory: → ② train AI models
③ Scientific sims
④ play games

Tensor core: 4th gen - supports FP8 and BF16
(Improved perf for AI/HPC workloads)^{dtypes}

Hopper Streaming Multiprocessor :-

Building block of H100GPU.

- 128 FP32 core, 64 INT32 and 4 Tensor core

NVLink is High Speed interconnect to
enable fast data transfer b/w GPU and
other comps

→ Interpreting message forwarded by Horus

Specs: ① Peak FP64 tensor core perf
upto 160 TFlops (SXM5)
and 48 TFlops (PCIe)

- ② FP32 : upto 160 TFlops and 48 TFlops
- ③ FP16 : 1120 " 96 "
- ④ INT8 : 12000 TOPS and 3200 TOPS

Memory and BW :

- HBM2e mem - 80 GB/s BW of 2K GBPS (SXM5)
- L2 cache - 50 MB (SXM5) and 40 MB (PCIe) 1155 GBPS

A100 - Ampere arch.

H100 - Hopper

Streaming Multiprocessor: Building block of
GPU. Multiprocessor that executes multiple
threads concurrently (using a combination
of units, register and memory)

Hopper arch (SM) 128 FP — Scientific sims, DL and compute heavy
64 INT — data proc, encryption and work compression

Features: 4 Tensor core: Specialized for matrix mult.
and acc. → DL/AI workloads

TFLOPS: Teraflops - measure of GPUs peak
FP perf. # FP ops / sec

TOPS: Terlops → Peak INT perf (#INT op/sec)

Design of Dual Brain Teacher Agent

HBM2e \rightarrow High BW memory &
↳ Memory type in H100.

High BW and low latency

NVLink could've fastened data transfer
b/w GPU and other parts (say CPU)
— reduces overall training time

GRID K520 GPU

H100

- 1.3 TFLOPS
(peak single FP precision)
FP perf)
- 0.65 TFLOPS
(double precision FP perf)
- Mem BW 208GBPS
CUDA cores = 1536
- 160 TFlops (SXM5)
& 48 TFlops (PCIe)
- 130 TFLOPS (SXM5)
& 24 TFLOPS
- 2000GBPS (SXM) &
1155 GBPS (PCIe)
CUDA cores, 14592
(SXM5) and 5888
(PCIe)

123x Faster than GRID K520

Roadmap of GPU Arch:

Fermi \longrightarrow Kepler \longrightarrow Maxwell \longrightarrow Pascal

Hopper \leftarrow Ampere \leftarrow Turing \leftarrow Volta

$1080 Ti = 11.3 \text{ TFlops}$ (S)
[Pascal] $0.5 \text{ TFlops (FP64)}$ 160 TFlops (H100)
93.5k CUDA core 130 TFlops (SXM5)
484 GBPS Mem B/W

Tensor core, NVLink & PCIe Gen 5

Input

General chat

Docs search (tool)

Generator agent

Output

How to reliably connect LLM to DB?

- Function calling / structured o/p
NO Regex

class RouteQuery (Base Model):

• datasource: Literal['vectorstore',
'generalchat']

LCEI
gradee-chain = gradee prompt | llm.with
[o/p of Left → I/p of right] → structured output (needs
LangGraph Document)
Ans: structured responses

Estimator → probator → promotor → predictor

framework of open tools

MCP server - Defines connecting interface to tools / data

Old - write py fn. for each agent
New - server exposing tools - agent connects tools like USB device

LangGraph - Define state and flow!

chain: input → op (stateless)

they need memory of what happened in prev. steps

class GraphState (TypedDict):

question: str # OG i/p

generation: str # final ans

context: str # rawtxt in DOCS

Start → Node 1 → Gen

{"qn":

" - - - "}

Router

{"output": ?

How did you build a route?

→ used pydantic model RouteQuery
→ Structured output to ensure type safety
→ How did you manage state? LangGraph

why can't I shove NVIDIA documentation
into prompt?

- Cost
- Accuracy
- Latency ↗

→ Web Scraping vs Cleaning :

WebBase Loader powered by bs4

(Raw HTML - full of noise - HTML
tags)

RAG ace → Data cleaning - X

Strips tags and leaves just text content

Chunking :- RecursiveTextCharacterSplitter

↳ chunksize : 1000

chunk_overlap : 100

[Preserves semantic coherence]

↳ splits by paragraphs first

↳ then sentences

↳ if not, words

Embeddings : NVIDIA nv-embedqa - e5 - v5

Trained specifically for QA pairs
hence performs better on technical queries
than OAI embeddings

Design overview:

Tutor for SA - NVIDIA role

to know more abt products
of NVIDIA

- that should teach me and
help generate stuff in resume
- test me
- contain study mat. to
official NVIDIA does

Multagent system using LangChain

@ Router - Llama 3.1 - 70B to intent
classify

② Tutor - searches SOT Doc pg's
and stores in RAG

Retrieval : (used E5 embedding & search)

+ FAISS
Reasoning : Sequential thinking loop (Analyze, Plan
and Teach) - GT before gen. tokens

③ Contextual Brain / Generic Tutor -
Ingest PDF to vector index (Ingest skills
dynamic ah)

why LangGraph

- state machine (NOT DAG!)

Retrieval + thinking in parallel

[Streamlit's model visualizes it for debug]

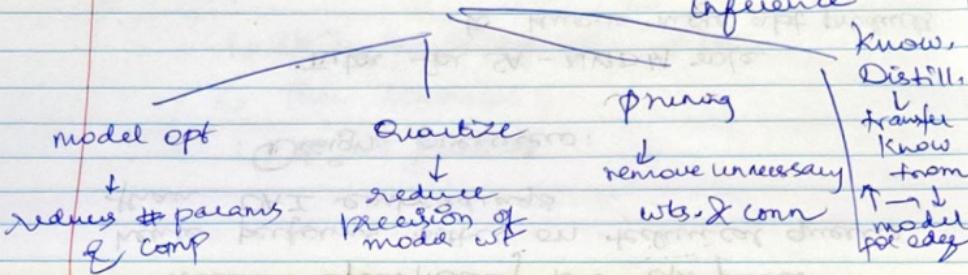
NVIDIA AI endpoints - connects to NIM
[NVIDIA Inference Microservices]

ChatNVIDIA → points to build.nvidia.com library

cluster of NIM containers

running on H100s

Tensor RT - SDK to allow developers to optimize for spec CPU arch, deploy on GPU for efficient inference



Agent

→ agent's train of thought
Memory [semantic sim, recently, importance, appn. spec metrics]

user — Agent core — Planning module

diag.

Tools

Core logic and behavioral char.

[Personas of agent]

ReAct, Reflexion, CoT, GOT : Reflection / critic

[Evidence based prompting

Framework]

Nitrus DB → better on perf. and scale on GPU arch.

→ expose docker to all GPUs

→ expose API key

→ NIM downloads + store model wt.

→ run as local user

to confirm service is ready --

curl →

curl localhost:8000/v1/health

/ready

Send a query in CLI:

Curl -X 'POST' 'http://127.0.0.1:8000/v1/chat/completions'

Request

```
-H 'accept: application/json'  
-H 'Content-Type: application/json'  
-d 'model':  
'messages': [ { 'role': ' ', 'content': ' ' } ]  
'max_tokens': 100
```

Nemo agent toolkit

- Framework agnostic
- Profiling / observability
- Eval. System

functions :

current_dt :

- type:

getting_started :

- type:

elms:

- prefix:

nim_hm:

- type: (nim)

model_name:

temp:

workflow :

- type :

o elm-name :

not_names:

@ register-function → decorator to inform
a toolkit that a fn. should be accessible by
name when referenced.

↳ specify config-type

Params of
decorated fn

} FunctionBaseConfig , builder

↳ dynamically query +
get other workflow fn.

NAT relies on yield return

- nat run --config-file
-- input "

async with load_workflow(config-file) as workflow:
async def single_call(input_st):
 async with workflow.run() as runner:

return await runner.
 result()

(yield)

async for callable in get_callable_for_workflow():

```
Langchain-nvidia-ai-endpoints import
```

```
ChatNVIDIA
```

```
llm = ChatNVIDIA(llm='meta/...')
```

```
tools = [TawnySearch(max_results=1)]
```

```
llm_with_tools = llm.bind(tools=tools)
```

```
llm_with_tools.invoke( ) . tool_calls
```

```
{ t.name: t for t in tools }
```

```
tool_call = ai_response.tool_calls[0]
```

```
- {t.name: t for t in tools}
```

```
'agen'
```

```
add_edge('---start---', 'agent')
```

```
add_conditional_edges('agent',
```

```
should_continue)
```

```
workflow.add_edge('tools', ).
```

```
app = workflow.compile('agent')
```

```
def should_continue(state):
```

```
messages = state['messages']
```

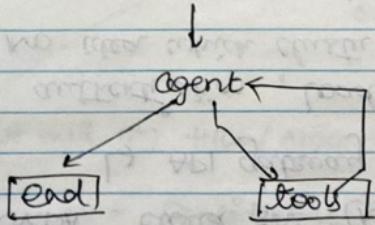
```
last_message = messages[-1]
```

```
if last_message.tool_calls:
```

```
return 'tools'
```

```
return '--end--'
```

```
Start
```



```
for chunk in appo.stream(
```

```
{ 'messages': [ ('human': 'what's the weather in SF?') ] },
```

```
stream_mode='values'),
```

```
chunk['messages'][-1].pretty_print()
```

```
chunk['message'][-1].pretty_print()
```

```
NIM handles model_selection, optimization  
+ backend selection and scaling automatically
```

```
Basic docker commands
```

```
* Optimizes LLM deployment
```

- OpenAI compatible REST APIs
- Deploy models from HF Hub
- Deploy models from local storage

Inference backends :

- a) TensorRT-LLM (high perf)
- b) vLLM (versatile deployment)

Container runtime → Docker → NVIDIA container runtime

NIM } Request lifecycle :
API }

Client → NVIDIA cloud → Orchestrator → NIM
(K8s) container

NVIDIA cloud first (NVCF)
↳ API gateway (severless) } GPU
authentication, load balancing & routing
No idea which cluster you're running

Orchestrator (Kubernetes) → DAX cloud
↳ routed to specific pod on NVIDIA NIM container
Auto scale → spins more pods & GPUs (K8s)

- NIM Container - highly optimized server stack

① API server

FastAPI/uvicorn
(OAI JSON req to tensors)

② Triton Inference service

(manages queue -
100 rays in, Triton batches them
(Dynamic batching))

③ Tensor RT-LLM

↳ Doesn't run raw PyTorch
code - runs optimized plan
for H100s

Triton on H100/A100 GPUs

Why local deployment?

- Confidential data - medical data
- latency - RTT to web
- cost @ scale → Rent / token M1
Rent GPU - flat rate

ChatNVIDIA : langchain_nvidia-ai-endpoints

integrate api.nvidia.com → localhost:8000

Request hits NVCF API gateway → routed to K8s pod

LIFECYCLE

executes
optimized
kernels

passes to
Tensor RT
engine where
mat-mul happens
(Optimized for H100s)

Triton Inf service
batches req

PreAct → more LLM calls

- prompt sensitivity and tuning overhead
- possible risk of hallucination
- lack of helium
- complex long chain

Toolcall: user query → fn-match → tool exec

→ Needs an LLM that supports
tool call

↓
resp.

not run —

serve = server that listens to incoming
eval = accuracy ↗ req.

functions: → wikipedia search:

llms: -type:
workflow: max-results:

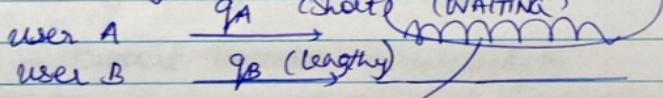
num_llms:
-type = num
model-name:
temperature,

recall True
recall get
recall get
recall get
response get
response get
user_name: user_name []

Tensor RT-LLM: not just a runner. It's a
complier.

↳ takes in llama3 and rewrites
math to perfectly match TH00 APU.

Standard batching



BLOCKING

(A) Soln :- TRI-LLM evicts user A's req
the moment it's done and slots in user C
even while user B is generating

Max^m. GPU util - we don't pay
idle compute → Dynamic / continuous
batching

(B) Paged Attention (KV cache optimization)

— LLM needs to remember prev. Tokens
[KV cache]

— Standard mem-alloc reserves huge blocks of
VRAM in case some gets long
(wastes 30-50% GPU mem)
↳ FRAGMENTATION.

Solution architect

[Pre-sales]

Architect - graph pieces together and bring them into a working sol'n.

work in engineering to fix bugs

Mission - support broader community to adopt AI technology

SW engg. comp than HW manufacturer

- = multiple triton inference server (customers)
- = work in prod mgmt to resolve → Engg. calls
- = Support proposal dev. (sizing, scaling etc.)
prompt engg.

NVIDIA product knowledge:

- NIM agents blueprint (reference workflow)
- Model served as NIM container
- NeMo guardrails to ensure model doesn't o/p proprietary code to public server
 - Deploy on triton inference server running on H100s to leverage transformer for FP8

How would you explain about NIM to someone who knows Docker?

NIM - repackaged - TensorRT, Triton, model wt. compiled for GPU arch,

docters own never io / nim / llama3
All you need is this - no need to install anything

The Nemo Framework (vs) Guardrails :-

Creation phase (Training and fine-tuning)

Runtime phase (Safety layer b/w user and LLM to enforce topic, safety and security)

MCP protocol:-

- standardizes AI agents to connect to data sources
- Instead of 50 API connectors, you build one MCP server - solves n^2 integration prob

Optimize a reasoning agent :

- FP8 (H100s) or INT4 quantization using TensorRT LLM to reduce mem. bw, pressure
- Optimize system prompt to minimize reasoning steps
- Speculative decoding, small model to predict tokens and big model verify them

RAG system that does real-time update :-

- Real-time vector DB: Milvus / Weaviate

Standard RAG: slow to index

Hybrid search: BM25 + vector search

↓
for most recent logs (kw search)

Real-time req. is met w/o waiting for embedding job to finish

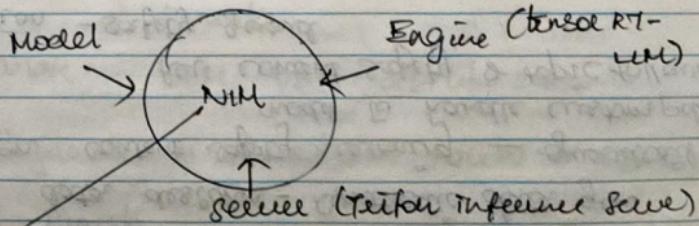
- Do we build MCP servers for users?

- AI Blueprints / custom arch. from scratch

Why NVDA?

→ NIM, NVLink, H100, TensorRT LLM,

NVDA → OS for AI.



[LLM in a box] → You can run it anywhere!

OpenAI → Run only on OAI server.

① Portability ② Data privacy (Deploy NIM within company's firewall) → e.g. banking) ③ Customization

You can't tell chatgpt to use FPS to save money

Portability

- Sovereignty
- customization

NEMO data designer - synthetic data gen
Nemeton content safety reasoning - guardrail model to handle custom policy for content safety & topic following
Nemeton - safety guard

Reducing KV cache bottleneck:

- grows ∞ prompt length and strains GPU memory.
- limited GPU mem.
 - (a) tradeoff in scalability
 - (b) reduced perf
 - (c) exp. recomputes
 - (d) more GPU

NVDA Dynamo 8 offload KV cache from GPU mem to CPU RAM / SSD / NV storage using low-latency NXL transfer lib

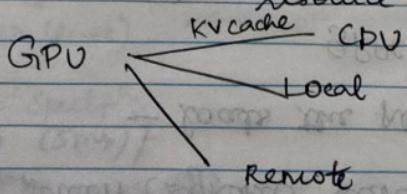
Dynamo's KV Block mgr:

- (Routed)
- separates mem mgmt from model engines, standardize storage access,
 - simplifies integration for scalability

offload / reload KV cache @ extremely high throughput, min. latency and maintain perf

Dynamo \leftrightarrow LM Cache (flexi cache mgmt) [Backend - CPU, SSD, Redis, GDS, InfiniStor]

KV cache offloading: most effective in big-context, high concurrency and resource-constrained env.



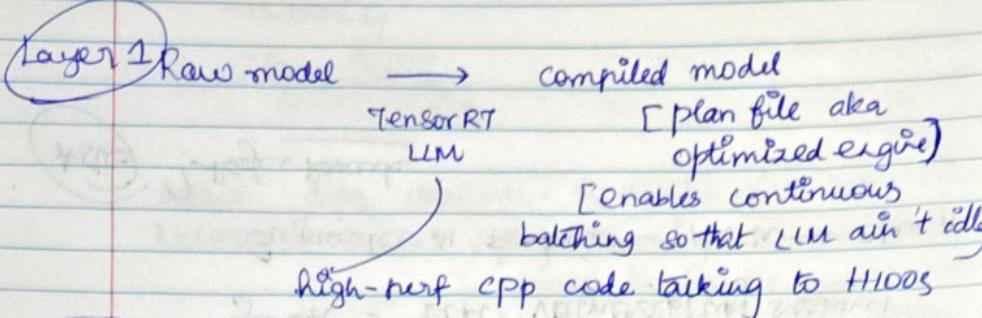
garak - LLM vulnerability scanner

Gigaccale AI factories - 100W \rightarrow OpenAI

Astra

? Raghu Ramaswamy

OpenAI Triton : DSL & compiler to write high-perf custom GPU using Python



Layer 2 Triton (Dynamic Triton) Node / Server

- loads the plan file → opens http/ gRPC ports

App. server (Nuvia, Unicore) → Req comes in ; Triton manages priority queue uses TRTLLM to do optimized math

Layer 3 cluster / Fabric ; Dynamic manages fleet of instances running Triton
[Basically a router]

→ Instance A has the memory (KV cache), so route there

NXL - moves data b/w nodes

Load balance + tg (VPC NW layer)

Instead of building 1 and 2 layer, use NMs

Diff b/w continuous batching and dynamic batching :

(Triton) Request-level packing [Server level]

waits for specific window (5ms)/
Spec # of reqs (16), once met, sends all of them as a batch

10 words → short time

vs 500 words ↑↑↑ time

Tput / latency ?

Continuous (TRTLLM)

↳ token-level packing
One req over

↳ insert next req.

GPU never idles
(2x/4x >> dyn batch)

Build a deep research agent

* LLM → RAG → Deep Research agent
[summarize, reflect]

Agentic capabilities in RAG

- Cons -
in general - no access to on-prem data
- unavailable for work use
- hard to understand blackbox

Deep research

Agentic AI

- Attention model +
Agentic reasoning + toolcall

NB 1 { NVIDIA apikey on Nae

Tawily key for web search toolcall

NB 2 { NAE deployment

Big model + more data → ↑ perf

NAE - works better P100 vs RTX

pre-training post-training test-time

compute

↑ PCT Scaling → enhance acc. on domain spec data

QA / Inst. following

→ open book exam + add. time to COT and do reasoning steps

~~Reasoning models for agentic AI app?~~

NB 3 Deploy RAG!

* Ingestion + retrieval / reranking + generation

NB 4 Deploy researcher.

Strawberry!

/no-think → no thinking

import OpenAI

from openai import OpenAI

Client = OpenAI()

base_llm =

base_llm =

api_key =

client.chat.completions.create()

model =

messages = [{ 'role': ' ', 'content': ' ' }]

see chunk in ...
if chunk.choices[0].delta.content
is not None

chunk.choices[0].delta.content is not None?
point (chunk.choices[0].delta.content, end=
✓
accesses 1st and only completion
choice
new txt fragment of this chunk
- prevents new-lines

Power of reasoning Response quality → transparency → self correction → resp.length

understand docker flags:

remove cont when stop
docker run -it -d -v /path/to/local/nvidia.conf:/etc/docker/runtime.conf
Container NVDA runtime file
GPU access
Cache model wt. to avoid redownload
- name local-nim
- runtime=nvidia
- gpus all → allocate all GPUs
- p 8000:88000 part for API access
- e → pass API key file model auth
- shm-size
↳ local shared mem for multi-GPU inf.

docker logs -f local-nim

name
integrations, api, nvidia, com N1 / chat/completions
apikey not needed for local deploy.
0.0.0.0:8000/n1
docker stop local-nim.

NEMO retriever embedding NIM:

text → vector emb. for sem. search
Re-ranking NIM & Imp. retrieval acc.
by re-ordering search results

NEMO reasoning NIM: Reasoning model

Milvus Vector DB: Stores and searches

document emb. @ scale

Doc. ingestion svc. ?

RAE query service :

1) Extraction phase - doc. ingestion + processing
vector emb. - capture semantic meaning of content in high dim space
(2 close / similar vectors)

2) Retrieval phase - query proc and resp. generation
Search on semantic meaning than simple kw match.

Pipeline :- DB user authentication → RAE → Doc

⑤ Deploy NLU & NVIDIA endpoint:

L A Y E R S		
Infra	Proc	RAG
- MILVUS standalone	- NV-Ingest runtime (ML doc. extr.)	- RAG Service (Query/retrieval)
- MinIO (obj store)	- Ingestor service (Doc. upload and Ingestion)	- RAG playground (Web UI)
- etcd (metadata store)		
- Redis (cache for logs / chat hist.)		

NIM local deployment -

```
chatNVIDIA (base_net=''),  
model=  
temperature=
```

```
NVIDIAEmbeddings (model = ' ',  
docker run -d --name meta-clamav3-8b-  
metainf --gpu=  
-e NAC-API-Key -v "/LOCAL-NIAC  
CACHE":/opt/niac/cache  
-u $(id -u) -p 3000:2000  
nacx.io/niac/meta/clamav3-8b-  
metainf)
```

documents = []
for url in urls:
 doc = html_doc_loader(url) → // uses BS 4.
 documents.append(doc)

```
text_splitter = RecursiveTextCharacterSplitter  
(chunk_size =  
chunk_overlap =  
length_fn = len  
)
```

```
docs = text_splitter.create_documents(documents)  
metadata = [doc.metadata] → [document, metadata]  
index_docs()
```

```
model = NVIDIAEmbeddings(model = ' ',  
truncate = 'END')
```

for doc in documents:

```
    text = splitter.split_text(doc)  
    page_content = page_content  
    metadata = [doc.metadata]
```

```
    if os.path.exists('...'):  
        update = FAISS.load_local('...')  
        update.add_texts(texts)  
        metadata = metadata
```

load_local
update.save_local()

```
FAISS.from_texts(texts, embeddings, metadata)  
docsearch.save_local()
```

OpenTec AI:

Intro to RAPIDS:

Complex Task → process

Data loading
+
ETL

Model training
+
Analysis

Model Inf
& Deployment

data acq.
" cleaning
! preproc

- visualization
- training
- eval

- model deployment
- " management

- [Rapids for CuDF]
- Apache Spark
- Nemo
Data curate
dataframe

- DL
- CuML (RAPIDS)
CuGraph
(graph
analytics)

- TRITON (Inference server)
- RAFT
(vector search
ML)

Challenges →

`df[[a,b]].value_counts().groupby(a)`

indexes 'a': .head().
agg('max'). sort_index()
rest_index()

lambda x: x[

* - agg (lambda n: n.idxmax(1))
- aggf (lambda n: n.idxmax(1))

```
df[[a,b]].value_counts().reset_index()  
    • agg('count'; 'max',  
          'set b: first')  
    • set_index()
```

`df['Issue Date'] = df['Issue Date'].dt.weekday.map()`

df['o'] = df['o'].apply(lambda x: pd.datetime.strptime(x, '%Y-%m-%d'))

df_group = df.groupby(['issue_weekday']).sum()['n']
- count(), sort_values()

%load_ext cudf.pandas

%. cudf. pandas. profile

"line-peop-le

s and DataFrame replaced by proxy obj. that dispatch ops. to CuDF when possible

Culp doesn't support

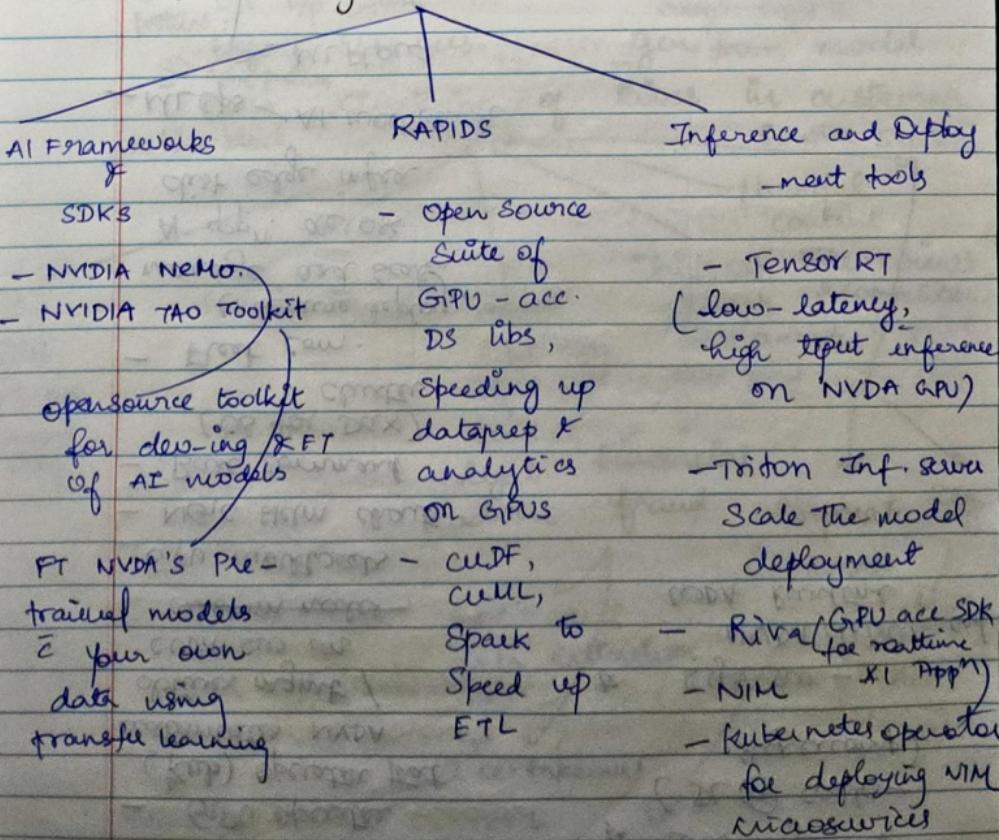
axis = 1 for `cant()`
method

~~df → df.groupby('').size().sort_index()~~

```
df.where(df[''] == '').groupby('').size
```

$$\Rightarrow \left(\frac{pc}{c} \times 100 \right) \rightarrow \text{pickup fraction}$$

NVIDIA's genAI toolchain :



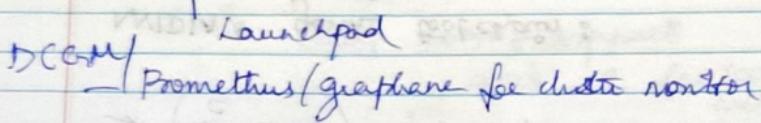
Orchestration and mgmt

- GPU operator
(Kub) operator that automates NVDA
 - driver mgmt / CUDA SW on custom nodes
 - GPU workloads
- NGC Helm charts
- Base command (OS for Dax) cluster
- Fleet com.
 - (Secure deploy, manage and scale)
 - A appn. across dist edge infra

Deployment and Infra layer

- * Dax system (DL @ scale supercomp)

- * Sys SW - NVDA driver and CUDA runtime



Nemo guardrails:

rules:

input: → Flows, blocks, patterns } config patterns

flows:

- detect-pii: Operates on i/p & o/p

blocks:

- detect-pii: intercepts unsafe req type regex on i/p side

patterns:
before LLM sees

- SSN pattern

- ph.no.

Clustering of Risks in customer service system

action: block

SSN pattern: [0-9]{3}-[0-9]{2}
ph.no.: [0-9]{10} - [0-9]{4}

Sanitize/block unsafe req. before gen from model

PII

Security bypass

legal

liability

- draft

lawsuit

- legal

advice

Jailbreak attempts

Inst override, Negotiate

Role play

O/p.

Financial

Fraud

Corporate Reg.

espionage

Violations

(HIPAA, GDPR)

credit card
data extraction

Acq. takeover!

- trade secrets

- pricing strategy

- Structured metadata for each block

Devs and compliance officers can see why it was blocked

`res = app.generate()`

`(res.explanations) → { 'blocked': true, 'reason': ... }`

- Compliance audits
- debugging
- risk dashboards

Longchain equivalent

- PII Middleware

- ContentFilter Middleware

- HumanInTheLoop Middleware

- SafetyGuardRail Middleware

after agent

Problems in AgentExecutor

- Blackbox

- state mgmt (can timetravel to any step)

Linear Loop

use this → `io.StringIO(generated_table)`
to convert
string to a table

Context free grammar

eg `CFG.fromstring(""`

`jsonschema.validate(instance=instance, schema=schema)`
`exceptions.ValidationException as e:`

DAPT: Dask, RAPIDS and cuDN

Download +
extract text → Lang det → Doc-level
→ Qual Filt

Content

Filt /
analysis

↓
Dataset

const
Blend +

Shuffle
Save Data
for custom
Tokenization

Streaming web + file queue

unified rep. of all data
text + metadata

metrics for pre/post training

Dedup, domain correction, domain validation

→ BOOST Acc.

* data diversity - MEST for diff. format, sources and dist. Ein domain

under-represented edge cases - present to prevent model bias

Synthetic generation of data (augmentation)
Creates diverse and balanced datasets

- Continuous data curation pipelines help maintain relevance

Arxiv, wiki, git →

Fuzzy semantic dedup

word count / line count heuristic
unide

save / shuffle to send to tokenizers

Downloader

Iterator

Executor

Run - curation pipeline → text + pdf

Sequential class - add built-in + custom filter

↳ read thru jsonfile

DocumentDataset ()

.read - json

ds_text = pass original dataset text
ds_code = " " code

tokenizer ready for model - training

curated

curation steps - text = Sequential (clean and unify, ScoreFilter (TextLineCountFilter), text_field = 'file-type-count', filter_text)

clean_type = 'bool'),

Domain - adaptive PT.

- learnt alt data creation

Sequential class

[Document Dataset]

jsonl → df. [Dask/Pay]



↳ unicode, PTI, Score, wordcount.

Modify() → filtering/manip

Custom tokenization :

- compare vocals of Oa tokenizer & newly trained domain spec tokenizer
- download llama2 embedding model & tokenizer

- Expand Oa tokenizer & new tokens to get extended tokenizer

- apply it to domain spec dataset, analyze user freq of new tokens and choose 98% cum freq

- init emb. of new tokens using Oa tokenizer → tokenize new token using Oa emb & o/p ots are found by avg

- merge new ots to Oa table to get DAPT.

SFT:

* Enables PT model to specialize in a gndomain

Style

- train on labeled data
- preserve broad knowledge acq. in PT.

* updates large portion /all of model w/

* Improved task spec. perf

↳ Needs relatively less data to train

↳ Needs task spec. i/p → o/p dataset

* Few hrs of modest HW.

Single SFT → multiple use - cases!

→ perf., cost & ctrl.

Model customization:

Prompt learning

FM

SFT (teaches model to follow style)

Foundation model

RLHF (user pref to ranking)

use model as to how TJs used

Fine tuning + eval

model → start from DAPT → or inst FT models
tokenizer → domain adapted tokenizer
tokens → Data, tensor, pipeline

Loss fn → prod. loss \in optional regularization
(cross entropy) \propto div

Learn Rate → smaller than DAPT → unsupervised

Risk of catastrophic forgetting
Easy to overfit on small high Q
created SFT data

Overfit prevention → ~~task~~ early stopping \in cross validation
MoE arch.

Eval → HF : (RLHF) \rightarrow accurate fb.
→ slow & collect feedback
→ subjective bias

LLM as judge → quick
→ biased by training data
→ prompt sensitive
→ hallucinate on Dom Spec qn.

B → 50% Tput ↑
200 } than
inc. mem & compute bco.

BEST TPUT : SFT + LORA ✓✓✓

✓ SFT:
=

JSONL → bin / idx → for faster access
good for storing & structured data

and memory mgmt.

needed for high tput data loading!

- allows direct, randomized access to data samples

(↑↑↑ I/O ops. reduces mem. footprint
Max 1-to-1 w/ idl, avoid bottlenecks in
data proc (efficient pretraining proc))

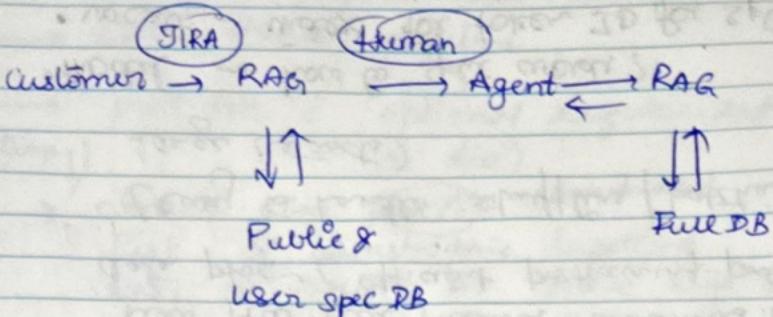
(easy to handle shuffling / batching of
large datasets)

• model - how to split words?

• vocab - index for token ID for split words

Sizing LLM Inference Systems

Learning # of GPUs needed for LLM workflow



Q: How many GPUs are needed?

→ gather requirements:
* 3.5k words in
 500 words out
* NeMo 43B GPT
* TTFT: 3s (limit)
* max: 31 req/s

Inference sizing: 13DGX + 100 systems

- Tput: 2.4 req/s
- latency for 1st token: $2606ms = 7.2s$ (Prefill) [within limit specified]
- Inter-token latency:

Tput, latency, TCO metrics

Prefill: time b/w pressing enter and first

%/p token (TTFT)

<20%

- Depends on # of i/p tokens

of end to end - compute bound

- populate KV cache for all tokens from query

[OR] prompt

{ Decoding - memory bound

- generate resp word by word (and token by token)

80%

Streaming

INFERENCE SIZING:



- matters to ppl how quick they get resp

util, max batch size

- low latency req - decreases available tput

Streaming
token - one @ a time
- TTFB *

Sequential
- full resp.
- E2E L

Qns. to ask for Inf. sizing :

- What model?
- avg # tokens in prompt to LM?
- avg # token in O/P
- Req/s? @ peak
- Latency limit - First/last token
- GPUs → what do you consider?

Prefill — Decode

~~compute KV cache in one shot (load weights @ one shot) — high work / load~~

autoregressive flows — one token @ a time

↳ load wt once — perform tiny
 $mat \times vect \rightarrow$ share in
cache

high mem bw use → compute to
load wt for each token

- Decoding is serial by nature

(mem BW)

- Balance Compute + memory bound

IFB : process prefill + decode in same batch

- nearly const. batch size for each token
↑ GPU util

- new req. exec quick as it waits for token
gen ONLY, and not completion of curr. reqs

Chunked context : separate prefill into chunks
batch one prefill chunk + multi-decoding
to attempt a balance
(implemented in TRTCLM)

- Irregular batch size introduce variability
when measuring latency

If concurrency = c , client sends c reqs
in parallel

→ Maintains conc = c

→ c outgoing concurrent reqs.

Max. batch size :

Benefits of TP:

- N times low mem. footprint / GPU
- N times TT mem b/w
- $\times N$ - compute resources / model
- Same precision, same acc.

TP2: comm. overhead b/w GPUs

→ Need low-latency GPU interconnect

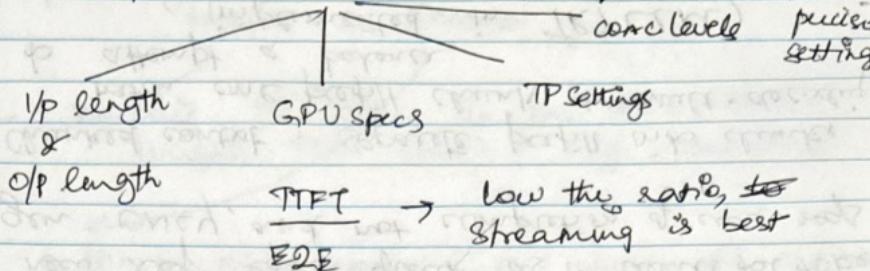
Reco. - NVLink enabled servers like

DGX, HGX

- pairwise connected PCIe cards like

H100 NVL

Parameters to tune for tput

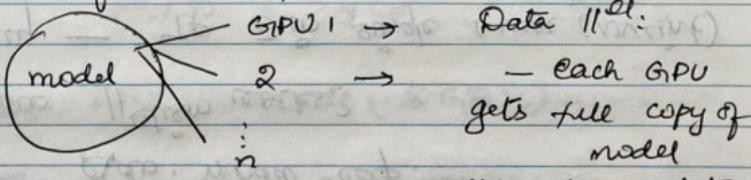


FP8 (Same Hoppe): — Reduced storage, inc tput & min. accuracy loss

- GenAI perf + NINM benchmarking guide to measure performance
- models with $> 13 \text{ B}$ params, use NVLink enabled systems
- Direct resp of LLM = stream
- Consumed by other LLM = seq.

cost and latency

Types of model parallelism:



— diff batch of data (for spec. gen.)

* works iff model is small enough to fit into a single GPU.

Goal: MM Tput.

Tensor Elision: Reduces latency for 1 req. Small computation per core

Model is too big for a GPU, split math opn.

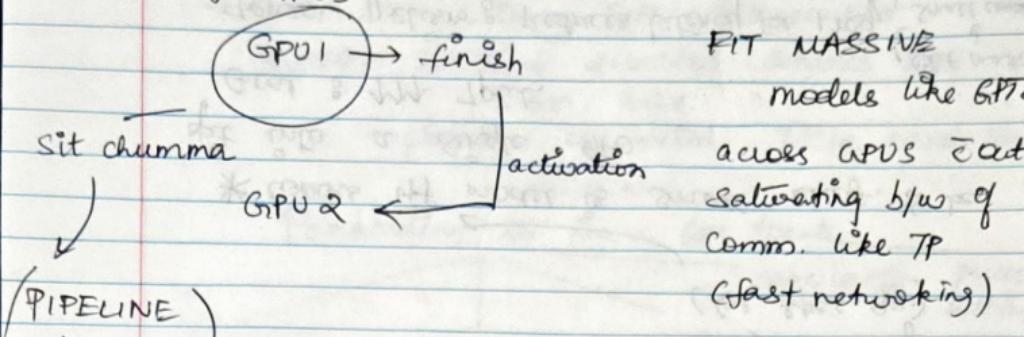
Split a tensor $[1024 \times 1024]$ by half

NYLink (GPU 1 calc. top half, GPU 2 calc. bottom) → combine partial results before moving to net layer

Pipeline Parallelism:

* If a model has 80 layers, we can't do TP. Talking b/w GPU is slow.

* Split the model layers [1-20, 21-40, ...] to diff GPUs.



Adv: model cap.

How PPI works?

- low latency — TP in single server (nvlink)
- model cap — PP access mult servers to fit big model
- DP across cluster to accommodate ~~many~~ users (concurrently)

For a 175B parameter model, start with

$$175 \text{B} \times 2 \text{ byte (FP16/BF16)} = 350 \text{ GB}$$

NVDA H100 → 80 GB HBM³ (High BW mem)

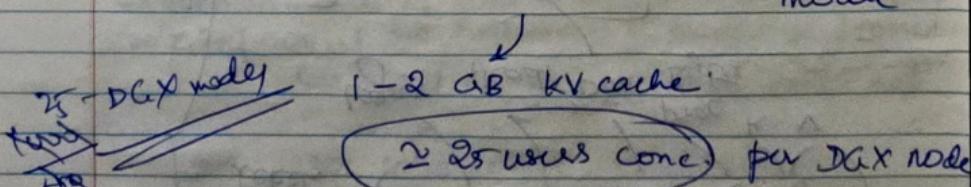
$$350 / 80 \geq > 4 \text{ GPU} = 5 \text{ GPUs}$$

~~DGX~~ DGX H100 node (8 GPUs)

$$\text{TP: } \rightarrow 350 / 8 = 43 \text{ GB / GPU}$$

FITS into H100. 37 GB free

KV cache
context window of 2048 - 175 GB model



need 40+ DGX nodes

Quantize model wt → 350 → 175 GB
only 2d GB for wt, 58 GB for KVcache

Evaluation and customization of LLMs:-

why eval LLMs :-

- boost user trust
- max^{m.} acc & reliability
- enhance safety and compliance
- prevent AI hallucination
- adapt to local lang

* Cloud (vs) on-prem

* eval error rate and ensure fairness

How eval helps?

- data driven insights
- error analysis

Trade-off b/w cost, latency and quali

validate ethical / bias eval over time

- Controlled eval - ✗

- Human-centered eval (HITL) - Human In The Loop

- Field eval ↳ relevance + language / tone
↳ A/B test, resp latency, error rate

✗ ✗ ✗

✗ user satisfaction, helpfulness

Controlled eval. — Text classification

language modeling
generation

QA

Machine Translation

Conc. context understand

- Dialogue
gen/creature
writing

Metric:

Quesn:

QA : Acc, prec, recall, F1 score
test

Lm metric: perplexity - how much the pred is surprising (low → better)

BLEU/ROUGE → mach trans

'n' gram overlap

~~flair~~ (summary score of ROUGE = human)

Academic /benchmark scales:

- GLUE
- Squad
- MMLU & F1SM8K
 - ↳ grade school math
 - 8K
 - math prob solving
- 8.5K open-ended math word prob.

③ qualitative / hybrid metric:

- Human judgements and LLM as a judge
clarity, coherence, factual correctness

(Pairwise comp, rank o/p, open-ended tasks)

* Task Spec. copy qual: coherence, relevance

* Model Robustness / Bias checks

↳ consistency, bias/fairness checks

Custom LLM eval

- uncover domain spec weakness/bias
- make improved model selection
- Justify cost benefit and resource allocation

generalized claim → specific claim

MMLU eval → measure multi-task language understanding

(evaluates gen knowledge & reasoning)

- MCQs (4/5 optn.)
- Factual acc, reasoning, dom spec und
- Diverse categories

ans ē high prob in logits - CORRECT ANS

Prompt → model → logits → highest

Why visualization matters?
Helps find trend

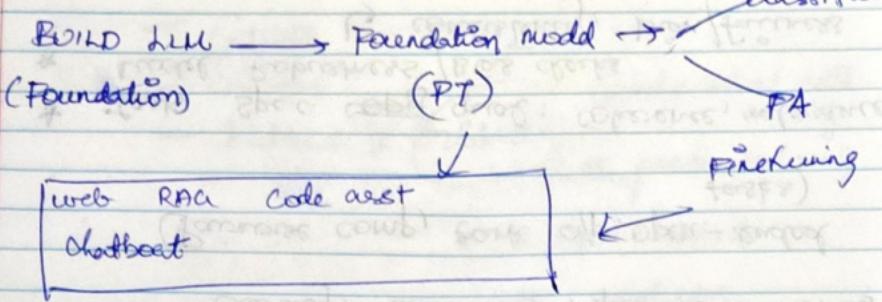
MLFlow! - Open-source platform

Log params, metrics, artifacts
and models

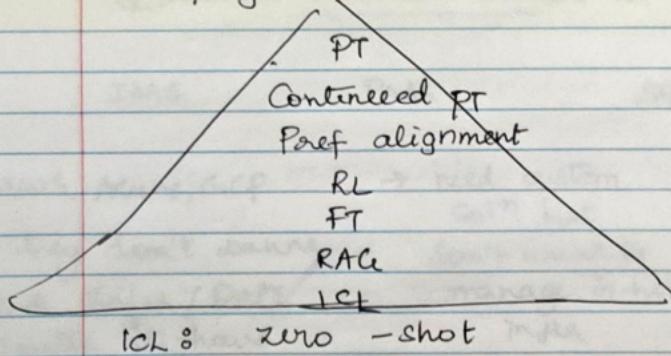
L-track eval across Run

Light customization of LLM :-

- Bridge domain gap
- Boost acc / reduce hallucination
- Naturalize tone/style
- enhance control / compliance



Adapting LLMs:



Few shot : prompt

SFT :- Train model params on labeled task data
- highest task acc.
- exp compute / large ckpt

Domain adaptive FT :-

Continued PT

Deep-dive into PEFT

- a) backbone frozen - Small # trainable param
- b) Techniques:
 - a) adapter - MLP module b/w layers
 - b) LORA - low rank decomp of wt. upd.
 - c) p-tuning - cont. prompt prepended
 - d) BITFit

↳ tune only

↳ key/val

Fast training, easy ckpt, simpler version ctrl

Business rounds - senior SA.

IAAS

PAAS

SAAS

Aws, Azure, GCP

- They don't wanna have Infra / Data centre in-house

pros:

- agility in provisioning
- on-demand scalability
- management of AI/ML capabilities w/o in-house expertise

→ need custom sol'n, but

don't want to manage in-house infra

Offer integration

svc. and build custom sol'n.

Enterprise Resource Planning (ERP)

Cust Rm.

Mgmt

(CRM)

Consolidate that data into central data warehouse

(Snowflake, Databricks)

- communicate to many stakeholders in an org
- high financial commitment
- legal/regulatory policy complexity