Planning the Al Architecture, Infrastructure & Operations

Pini Cohen Einat Shimoni



Copyright@STKI_2024 Do not remove source or attribution from any slide, graph or portion of graph

R (2)

5

This presentation will take you on a journey

Towards planning your organizations' AI architecture, infrastructure and operations





What are the needed architecture, operations and infrastructure to effectively deploy AI?



How should companies organize to leverage GenerativeAI models?



What are the main risks? How can organizations mitigate them? How will the EU AI act affect the industry?

Copyright@STKI_2024 Do not remove source or attribution from any slide, graph or portion of grap

Risks and

regulations









Current Data Analytics Architecture



STKI.INFO

Current Data Analytics Architecture Issues

Data Silos

Inflexibility

Poor Data Quality Inconsistency and Repetitive Scalir Security & Perfor Privacy operations-

e Scaling & Performance

High Costs



Which operational model do companies use today?



Existing data organizations look like this:



Issues:

- Slow TTM
- Dependency
- Focus on delivery
 (not discovery)



What are the needed architecture, operations and infrastructure to effectively deploy AI?







Copyright@STKI_2024 Do not remove source or attribution from any slide, graph or portion of graph





Which tools are needed for your modern data factory?

- Data Ingestion and processing-preparation : ETL, Streaming, CDC, code (python)
- Data stores Cloud \On premise (Vector Databases)
- Data virtualization tools (Logical Data)
- Data governance tools : QA, Security, compliance, lineage, data catalog, metadata management



- Self service BI tools visualization tools
- ML and AI tools
- Data pipeline & orchestration Data Ops. MLOps LLMOPS
- IDP Internal Development Platforms

 how will developers access analytics components





Copyright@STKI_2024 Do not remove source or attribution from any slide, graph or portion of graph

TKI.INFO

N





34% עובדים כיום 45% מתכננים להיכנס לעננים אנליטיים בקרוב



Copyright@STKI_2024 Do not remove source or attribution from any slide, graph or portion of graph

Ρ



Copyright@STKI_2024 Do not remove source or attribution from any slide, graph or portion of graph

TKI.INFO

Automation tools -MLOPS main features



MLOPS TOOLS LANDSCAPE

https://research.aimultiple.com/mlops-tools/

- Model building automation model selection, data preparation, feature engineering, model training
- Model registry / versioning
- Automation testing (& AB testing) and experiment management
- Model monitoring logging
- Model governance (risk, compliance)
- Model optimization (GPU / Cloud options)
- Continuous integration and deployment

The Operational management of LLMs in production environments. It encompasses the practices, techniques, and tools used to effectively deploy, monitor, and maintain LLMs:

- Data Collection and Preprocessing
- Model Selection and Training
- Evaluation
- Prompt Engineering
- Explanation and Interpretability
- Model Optimization
- Testing and Monitoring
- Deployment and Scaling
- Versioning and Updates
- Ethics and Governance

Automation tools – LLM Ops

The LLMOPS part from the landscape



https://research.aimultiple.com/mlops-tools/



Differences between MLOPS and LLMOPS

MLOPS	LLMOP
Mature	Less Mature
Smaller, label-based training data	Massive text corpuses requiring petabyte storage
Training from scratch or use transfer learning	Use pretrained Foundation Models with lots of customization options
N/A	Crafting prompts
Varies	Extremely large models with high computational needs
Sometimes important	Crucial
N/A	Block generated toxic text
Sometimes important	Very important
	MLOPS Mature Smaller, label-based training data Training from scratch or use transfer learning N/A Varies Sometimes important N/A Sometimes important



Copyright@STKI_2024 Do not remove source or attribution from any slide, graph or portion of graph

KI.INFO

Data Vault 2.0 - the middle road

Data Warehouse

Most of DW is based on Star Scheme

- Advantages:
 - Simplicity and Understandability:
 - Query Performance
- Disadvantages:
 - Data Redundancy
 - Lack of Flexibility

Data Vault 2.0

- Enables fast modification to the model
- Every data source is related to the existing model
- Fast uploads
- Primary Technical Keys are generated + source/time of loading

Data Lakes

Do it yourself. This makes everything faster. But:

- Data Quality and Governance
- Data Silos and Lack of Integration
- Security and Privacy Concerns
- Data Discoverability
- Performance and Scalability
- Lack of Data Ownership and Accountability
- Becoming "Data Swamp"







Copyright@STKI_2024 Do not remove source or attribution from any slide, graph or portion of graph

Ν

Israeli organizations are **just starting** to adopt a product approach for data & analytics



The Four Fundamental Team Topologies





Data Product Teams

A Data Product team is a "stream-aligned" team – several other teams including "platform", "enabling" and "complicated subsystem" teams make it easier for Data Product teams to do their work



Figure 1, Data Product Teams





How should companies organize to leverage GenerativeAI models?



Webinar בתאריך 5.2.24 נקיים על מונחי בסיס בעולם ה GenAI

Stki.info/events

GenAl is evolving fast – technology is not stable

GenAI model basic characteristics:

- Number of Parameters & number of layers
- Corpus Size the amount of text data the LLM was trained on measured in tokens



Foundation Model customization best practices are not DAT stable ALGRNTING

Prompt Engineering

- Iteratively refining prompts for better results.

Prompt Tuning (hard and soft tunning) - Optimizing LLM parameters based on specific prompts and outcomes.





Fine Tuning-Retraining the LLM on a small, task-specific dataset.

Parameter Tuning

RAG - Retrieval Augmented Generation

Retraining-

Training the LLM from scratch on a new, possibly large dataset

GenAI and your spending – It's complicated

Training a model:

- According to unverified information leaks, GPT-4 was trained on about 25,000 Nvidia A100 GPUs for 90–100 days
- Estimation is about 3,125 servers
 * 324K\$ per server were needed for total of 1B\$
- Estimation of training costs for GPT-4 was **around \$63 million**.



Soon we will need AlFinOps practitioners	Using a m Mistral Of summariz	odel (Inference) Den Source model ing text	
		A10 24GB GPU (1500 input + 100 output tokens)	A100 40GB GPU (1500 input + 100 output tokens)
	Best Latency	4.1 sec	2 sec
	Best Throughput (Max RPS without significantly hurting latency)	0.9	3.6
	Latency (at max RPS) (p50 / p90)	5.8 sec / 5.8 sec	4.7 sec / 4.7 sec
	Machine Type	AWS (g5.xlarge)	GCP (a2-highgpu-1g)
	Cost per hour (Spot)	\$0.30	\$1.45
	Cost per hour (On-demand)	\$1	\$3.50
AIFINOPS	Source: https://www.linkedin	.com/pulse/7-lartsim-30b-be	nchmarks-truefoundry-uj9bc/

GenAl issues:

- LLM's have tendency to generate plausible-sounding but false or nonsensical statements
- Hallucinations type:
 - Incorrect facts "was born in 2001"
 - Fabricated information "use this URL" "this company is doing ..."
 - Weird or Creepy Answers "I am in love with New-York"



https://builtin.com/artificial-intelligence/ai-hallucination



Example of hallucination

Question to LLM : give list of products for code generation? Answer: Kite: Provides intelligent code completions, documentation lookup, and on-the-fly code examples for Python, JavaScript, and other languages

Kite is saying farewell



By Adam Smith, Founder November 16, 2022

From 2014 to 2021, Kite was a startup using AI to help developers write code. We have stopped working on Kite, and are no longer supporting the Kite software.





What is the risk for enterprises using AI?

IP and copyrights

violations in the model recommendations , data trained by the model, Who has the IP rights for content created with AI help

<u>Security – business</u>

Bad/harmful Recommendations /business actions caused by AI hallucinations or by poisoned trained data



Security - technology

Malware inside Al output (especially code) Information is leaking (credentials, company secrets)

Privacy & regulatory

Toxic text generated , Data leakage of PII, Algorithmic discriminations



All the above in your supply chain

EU AI Act

Highest Fines (up to €40 million or 7% of global)



High-risk AI obligations effective: June 2025

Prohibitions effective: December 2024

Most provisions effective: Mid-2026

Formal adoption: Mid-2024



EU Artificial Intelligence Act: Risk levels





Which GenAI models/service currently comply with the AI UI ACT?

No existing (leading) models can claim full legal compliance to the EU AI Act yet*

"*No foundation model provider achieves a perfect score"



Stanford University Human-Centered Artificial Intelligence

Foundation Model Transparency Index Total Scores, 2023



Source: 2023 Foundation Model Transparency Index



GenAl is helping the big companies to become bigger

Foundation model training is very expensive

Customization models is better when you have more data





KIINFO

S





Defending the New LLM Frontier: End-to-End Security for the Generative AI Era

Tools / options to mitigate risks



- LLMOps capabilities (discrimination /bias engines / toxic text blocking, etc.)
- DLP when prompting or when using with API
- Compliance / security management tools for Opensource



- Use AI risk tools when available
- "Enterprise Edition" cloud solutions
- On premise installations



Recommendations for enterprises regarding AI risks

Evaluate their risks in the different vectors and set their policy to the different activities & domains

Enhance the security awareness program with AI related topics

Implement AI security tools when available

For most of the enterprises the "enterprise edition" offering might be good enough

Evaluate AI risks and solutions every 6 months

What are the future directions of the AI industry?

KI.INFO



6

What's

next?

A glimpse to the Al future





Beware of the "shiny object syndrome"

There is no substitute for:

A. Having a clear business caseB. Having trustworthy dataC. Continuous discovery





hank vou.

Good luck with planning your future **Al architecture & operations**



