# CIKLUM | SPEAKER'S CORNER
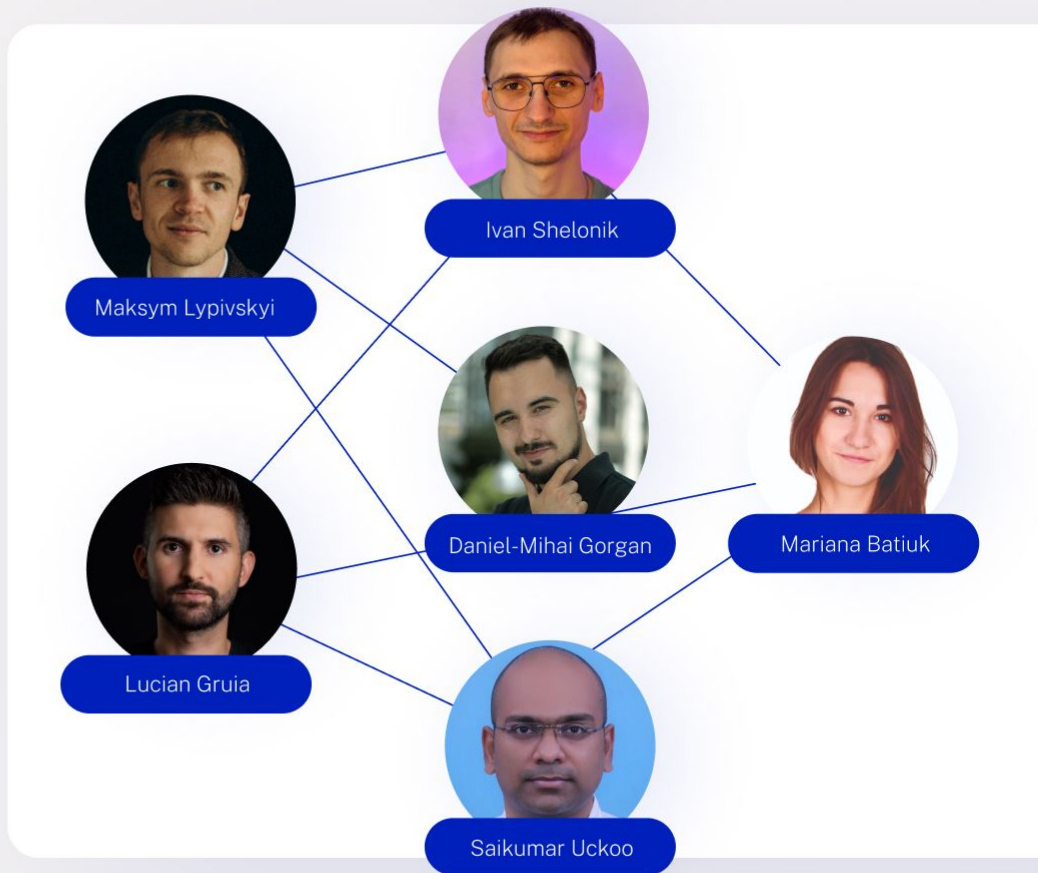
April, 16   16:30 (CET)   English

# Architecting Scalable AI RAG Systems:

## From Startup to Enterprise. A Live Coding Session



Maksym Lypivskyi

Ivan Shelonik

Daniel-Mihai Gorgan

Mariana Batiuk

Lucian Gruia

Saikumar Uckoo

# Experiences of tomorrow. Engineered together.

CIKLUM

We transform how people experience the business. All through next generation technology.

**What we do:**

Product Engineering | Intelligent Automation | Data & Analytics

**2002** founded

**4000+** professionals

**20+** offices

**300+** clients

## Leading companies choose us:

JUST EAT Takeaway.com | METRO MARKETS | eToro | ZURICH | Mercedes pay | Greensill KANTAR RETAIL | dacadoo

# Our Global Delivery Centres

Global Reach, Local Insight - Ciklum bridges the best in tech from the three key IT regions

CIKLUM

**Central & Eastern Europe**

🇧🇬 Bulgaria
🇨🇿 Czech Republic
🇵🇱 Poland
🇷🇴 Romania
🇸🇰 Slovakia
🇪🇸 Spain
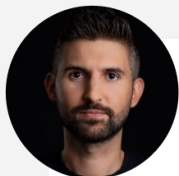🇺🇦 Ukraine
🇬🇧 United Kingdom

**Asia**

🇮🇳 India
🇵🇰 Pakistan

**LATAM**

🇦🇷 Argentina
🇺🇾 Uruguay

# Our speakers

## Lucian Gruia
### Principal AI Technology Lead

- AI Tech Lead with over 11 years of hands-on experience in Telecom, Fintech, and Aerospace. He specializes in AI, data integrity, fraud detection, system performance, architecting frameworks and solutions for real-time systems.

- Develops an AI upskill program for 300 engineers at Ciklum.

## Ivan Shelonik
### Expert Data Scientist

- Certified Professional Machine Learning Expert with 7 years of commercial experience in developing Machine Learning projects from the ground to delivery into the Cloud (AWS 5+YoE).

- Has worked and delivered primarily for customers from S&P 500

## Daniel-Mihai Gorgan
### Senior JS Developer

- Tech enthusiast specializing in Node.js, SQL/NoSQL and Cloud technologies with 5+ years of experience

- Hands-on experience in projects across outsourcing and product companies, contributing to the development of in-house products, smart chatbots, and voicebots by leveraging different AI technologies

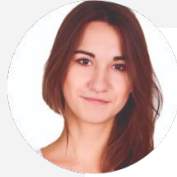# Our speakers

## Saikumar Uckoo
### Conversational AI expert

- Cloud Architect specialized in building, deploying, and maintaining AI solutions on Microsoft cloud platforms. Leads deliveries on platforms like Microsoft PVA, KoreAI, and custom GenAI solutions built on open-source tech.

## Mariana Batiuk
### Principal TCoE Lead

- Mariana leads the technical council on the QA maturity assessment, test strategy, pre-sales, new services development, initiatives, and quality engineering activities.

- Has proficient experience in QA Management, Agile Methodologies, Testing, Team Management, and Coaching.
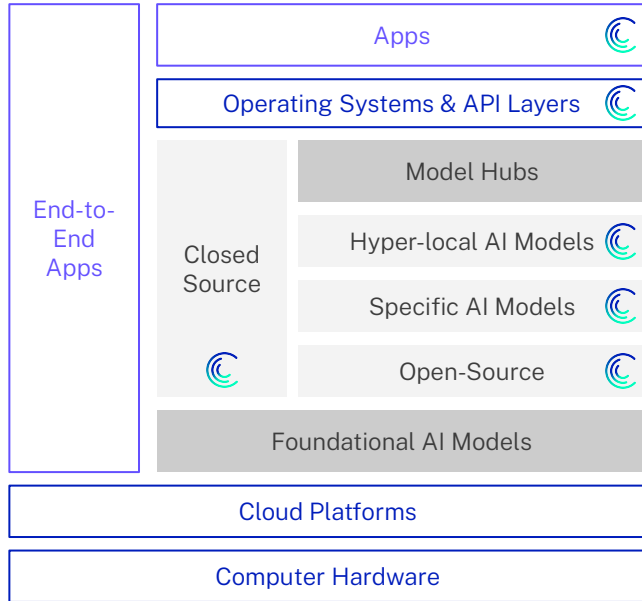
## Maksym Lypivskyi
### Global Head of Cloud Platforms

- Specializing in cloud computing architectures and generative AI applications, he focuses on creating, deploying, and optimizing cutting-edge solutions across global platforms.

- A mentor and community builder, he actively shares his insights on generative AI, cloud technologies, and leadership.

# Playing in all parts of the AI stack



## User Experiences & Engagement

**End-to-End Apps**

- Apps
- Operating Systems & API Layers

**Closed Source**
- Model Hubs
- Hyper-local AI Models
- Specific AI Models
- Open-Source

Foundational AI Models

Cloud Platforms

Computer Hardware

**Legend:** Applications | Models | Infrastructure | Where we work

## Emerging Stack Trends

### Applications
The rise of cloud-based generative AI and LLMs, accessible via **APIs** and **embedded** in other applications, will allow companies to use them as-is or **customize** with their data

### Fine-Tuning
The need for model fine-tuning will drive demand for a **diverse skill set,** such as software engineering, psychology, linguistics, etc.

### Foundation Models
The market will evolve and diversify with the emergence of **more pre-trained models,** offering **options** for size, transparency, versatility & performance balance

### Data
Mastery of **new and diverse data types** and volumes will be crucial for success, with **GenAI features in modern data platforms** facilitating **adoption at scale**

### Infrastructure
Essential for GenAI deployment, **cloud infrastructure** will help manage costs and **carbon emissions**, necessitating data center retrofitting and advancements in chipset architectures, **hardware** & algorithms

## Partners

- OpenAI
- Weights & Biases
- Humanloop
- HUGGING FACE
- Pinecone
- OpenAI
- appen
- redis
- snowflake
- databricks
- Azure
- aws
- Google Cloud
- NVIDIA.

# Agenda

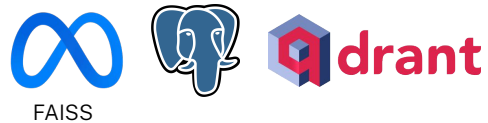| | | | |
|---|---|---|---|
| **01** | What is RAG | **05** | Build with Javascript |
| **02** | LLM Wrappers and Docker | **06** | Deploy RAG app in AWS |
| **03** | Build with Java | **07** | Deploy RAG app in Azure |
| **04** | Build with Python | **08** | Challenges in QA and more |

# Session's Tech map

## Programming languages

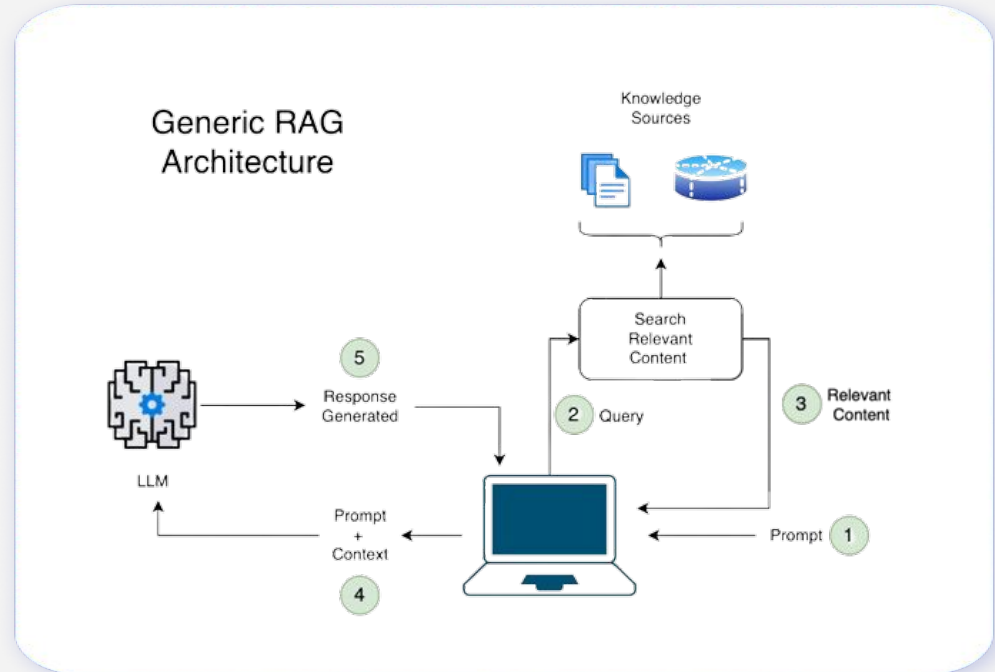## Databases

FAISS

## Infrastructure

# What is RAG

A RAG system essentially correlates a **user's prompt** with a relevant **data chunk.** It does this by identifying **the most semantically similar** chunk from the database.

This chunk then becomes **the context** for the prompt.

When **passed to the Large Language Model (LLM),** it enables the system to provide a relevant answer within the given context.



Generic RAG Architecture

Knowledge Sources

Search Relevant Content

5 Response Generated

2 Query

3 Relevant Content

LLM

Prompt + Context

4

Prompt 1

**9** 9

# LLM Wrapper

- Build with **Java**
- Deploy locally
- Integrate a 3rd party client

CIKLUM

Ivan Shelonik

Maksym Lypivskyi

Daniel-Mihai Gorgan

Mariana Batiuk

Lucian Gruia

Saikumar Uckoo

# Why do we need RAG?
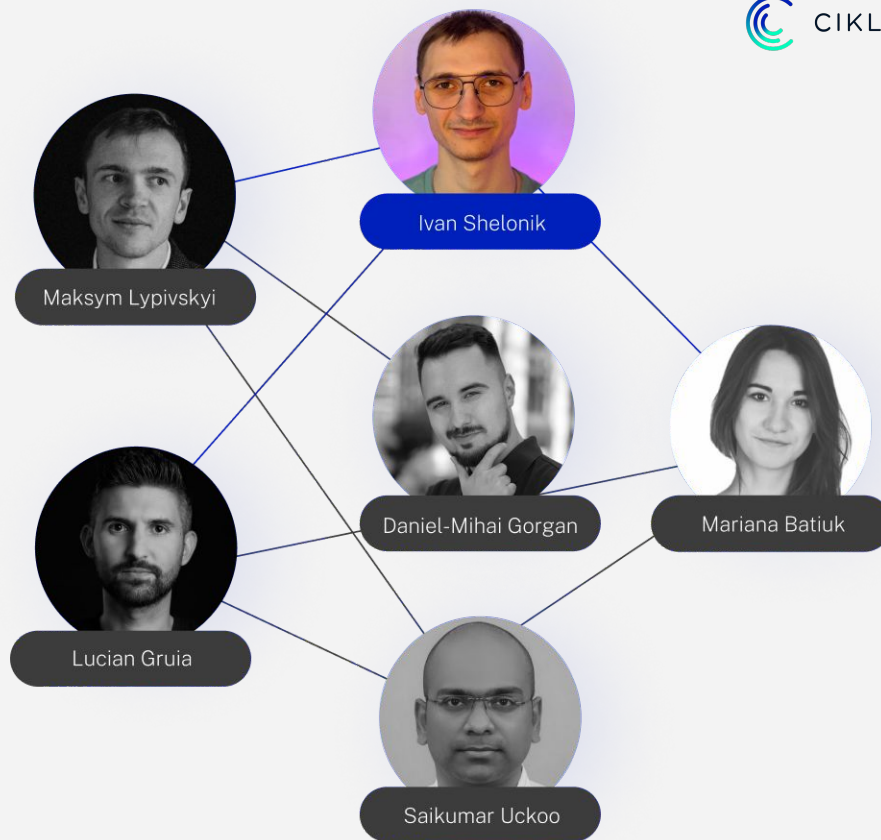
CIKLUM

- **Expands Knowledge Base**
  RAG accesses a vast external database, enriching its knowledge beyond initial training data

- **Improves Accuracy**
  Enhances response precision by integrating relevant, real-time information

- **Adaptable**
  Effectively handles novel and niche queries

- **Increases Efficiency**
  Streamlines information retrieval and generation process

- **Versatile Applications**
  Useful across various fields, from customer support to research



Source: What Is Retrieval-Augmented Generation, aka RAG?

# AWS

- Build with **Python**
- Build Docker images
- Semantic search with FAISS
- Deploy on AWS

CIKLUM

Ivan Shelonik

Maksym Lypivskyi

Daniel-Mihai Gorgan

Mariana Batiuk

Lucian Gruia

Saikumar Uckoo

# Data Chunking and LLMs

**LLMs** also have a limited capacity for context.
Just as humans **cannot digest unlimited context**, these models have a specific size limit for the content they can process.
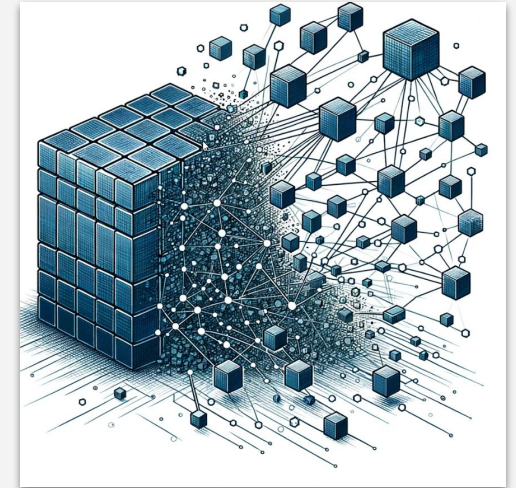
*So, what about situations involving very large amounts of data?*
Consider a specific use case, such as a book. It's too large to pass the entire book as **the context** for the current prompt, so it **needs to be divided** before being stored in the database.

This process is known as **data chunking**.

Types of Data chunking (by size):
- Fixed-size
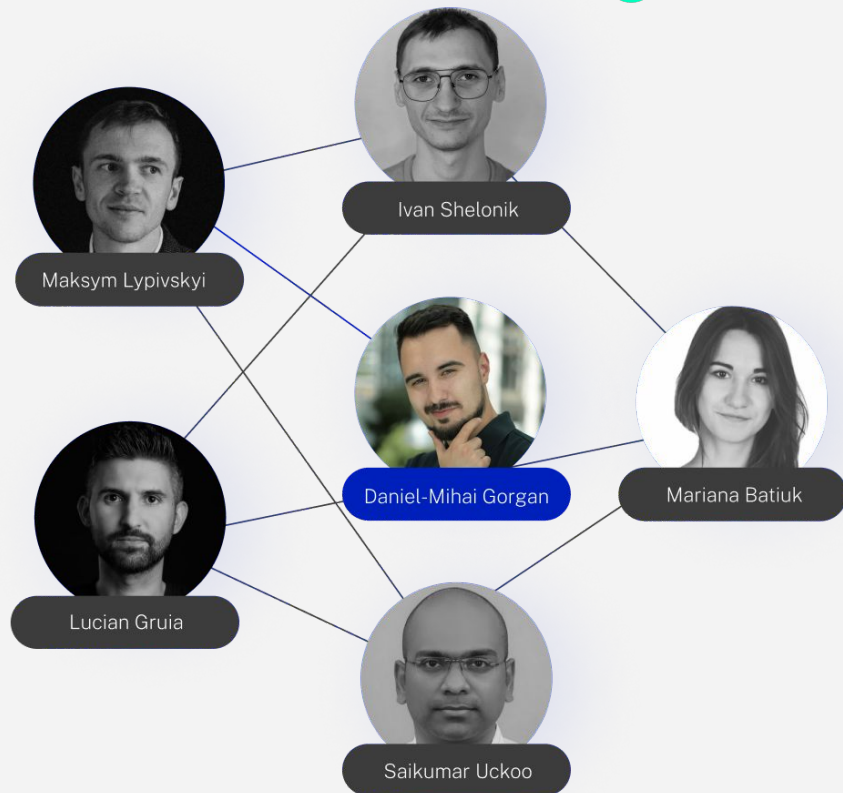- Variable Chunking
- Semantic Chunking



Generated with DALL·E 3

CIKLUM

# JavaScript

- Build with **TypeScript**
- Semantic Search with Pg vector

Maksym Lypivskyi

Ivan Shelonik

Daniel-Mihai Gorgan

Mariana Batiuk

Lucian Gruia

Saikumar Uckoo

# Embeddings. Similarity

- **Embeddings**
  Numerical representations of concepts, in a high-dimensional space, capturing semantic meaning.

- **Similarity:**
  - Lexical: entities are alike in appearance
  - **Semantic**: entities are alike in meaning

- **In RAG we represent entities by describing them.**
  This is a form of knowledge representation.

  **Example: Mountain, River, Canal**

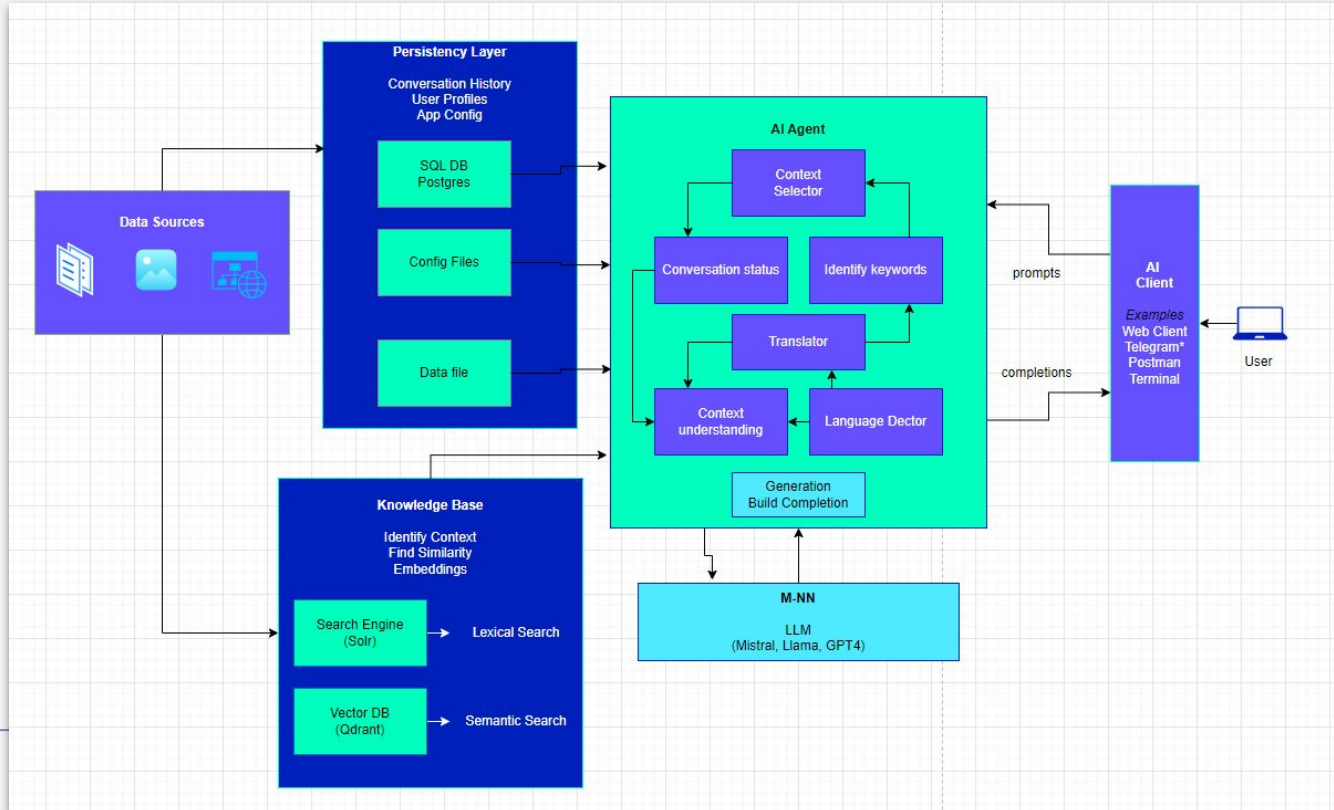| One hot encoding | | 2-Dimensional Space | |
|---|---|---|---|
| | | | [Natural vs Artificial, Mobility] |
| Mountain: | 1 | Mountain: | [-0.7,   -0.8] |
| River: | 2 | River: | [-0.3,   0.7] |
| Canal: | 3 | Canal: | [ 0.4,   0.5] |



Read more: Wikipedia - Cosine Similarity

# Azure

- Deploy on **Azure**
- Semantic Search with Qdrant
- Conversation history

CIKLUM

Maksym Lypivskyi

Ivan Shelonik

Daniel-Mihai Gorgan

Mariana Batiuk

Lucian Gruia

Saikumar Uckoo

# RAG Architecture

# Benefits of RAG

1. **Providing up-to-date and accurate responses**
   RAG ensures that the response of an LLM is not based solely on static, stale training data. Rather, the model uses up-to-date external data sources to provide responses.

2. **Reducing inaccurate responses, or hallucinations**
   By grounding the LLM model's output on relevant, external knowledge, RAG attempts to mitigate the risk of responding with incorrect or fabricated information (also known as hallucinations). Outputs can include citations of original sources, allowing human verification.

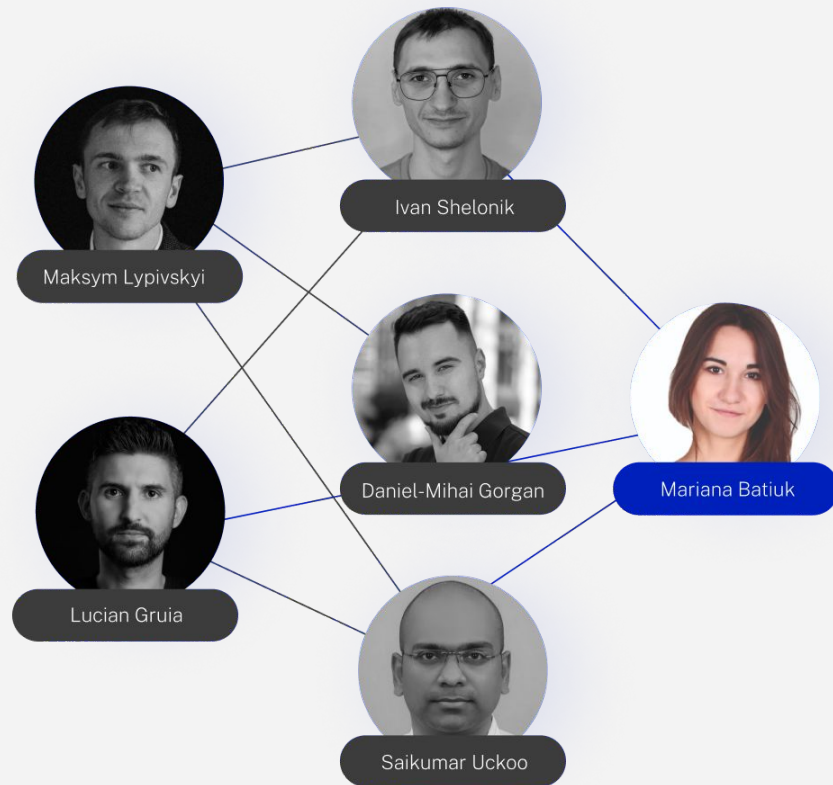3. **Providing domain-specific, relevant responses**
   Using RAG, the LLM will be able to provide contextually relevant responses tailored to an organization's proprietary or domain-specific data.

4. **Being efficient and cost-effective**
   Compared to other approaches to customizing LLMs with domain-specific data, RAG is simple and cost-effective. Organizations can deploy RAG without needing to customize the model. This is especially beneficial when models need to be updated frequently with new data.

# QA & Testing

- SW characteristics
- Top 5 **risks**
- Methods and tools
- Balanced **success** factors

CIKLUM

Maksym Lypivskyi
Ivan Shelonik
Daniel-Mihai Gorgan
Mariana Batiuk
Lucian Gruia
Saikumar Uckoo

# Software Characteristics
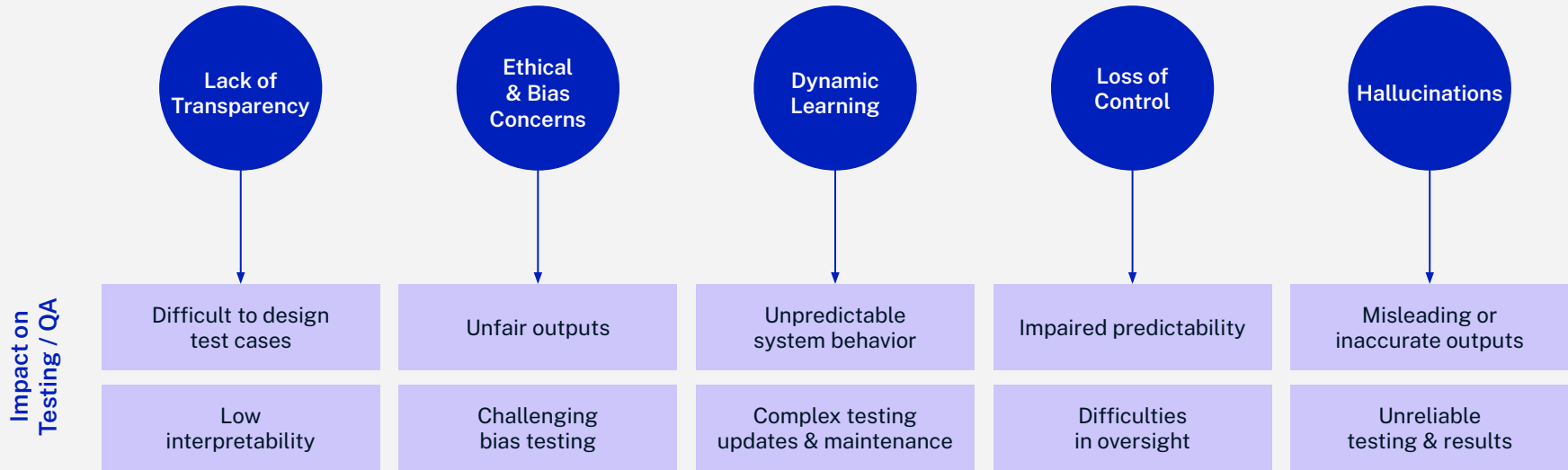
| ISO 25010 Product Quality Model | | | | | | | |
|---|---|---|---|---|---|---|---|
| Functional Suitability | Performance Efficiency | Compatibility | Usability | Reliability | Security | Maintainability | Portability |

Functional Testing → *What* does the system do?

Non-Functional Testing → *How* do the system do this?

| AI-specific Characteristics | | | | | | | |
|---|---|---|---|---|---|---|---|
| Flexibility & Adaptability | Autonomy | Evolution | Bias | Side-effects & Reward Hacking | Ethics | Transparency, Interpretability & Explainability | Safety |

# Top 5 current shortcomings and risks

CIKLUM

**Impact on Testing / QA**

| Lack of Transparency | Ethical & Bias Concerns | Dynamic Learning | Loss of Control | Hallucinations |
|---|---|---|---|---|
| Difficult to design test cases | Unfair outputs | Unpredictable system behavior | Impaired predictability | Misleading or inaccurate outputs |
| Low interpretability | Challenging bias testing | Complex testing updates & maintenance | Difficulties in oversight | Unreliable testing & results |

# Some essential elements

that should be considered when verifying AI systems

CIKLUM

KNOW
THE ALGORITHM

TEST
THE ALGORISM

BALANCED
SUCCESS
FACTORS

MAKE SURE TO HAVE
ENOUGH DATA

BRING THE
RIGHT PEOPLE

# Interact

- **Prompt** Engineering
- Fine-tuning

# The optimization flow

Context optimization

What the model needs to know

| RAG | All together |
|-----|--------------|
| Prompt engineering | Fine tuning |

How the model needs to act

LLM Optimization

# What is a good prompt

Act as an experienced Learning specialist. I need to improve my upselling skills. Prepare an educational program for me to improve that skills.  Program should be for 2 month with 4 hours effort per week.

Please provide answer with the next output:

Topic: Name

- blocks
- ...

Books:

Example:

Topic: Negotiation basics

- Win-win strategy
- Active listening strategy

Books: "Getting to Yes" by Roger Fisher and William Ury

Instruction
Context
Role
Formatting
Tone
Examples

# Prompt tactics

## * Shot Prompting

Zero
Add 2+2:

One
Add 3+3: 6
Add 2+2:

Few
Add 3+3: 6
Add 5+5: 10
Add 2+2:

## Model-guided prompting

Before answering, I want you to first ask for any extra information that helps you produce a better answer.

If you got no questions, please provide an answer instead.

## Self-evaluating prompting

Can this program be improved?

# Chain of thoughts

Virma has three bags, each of which fits five shirts. How many shirts can Virma fit in her bags?
Let's think step-by-step.

CIKLUM

**Prompt**

Question + "Let's think step by step"
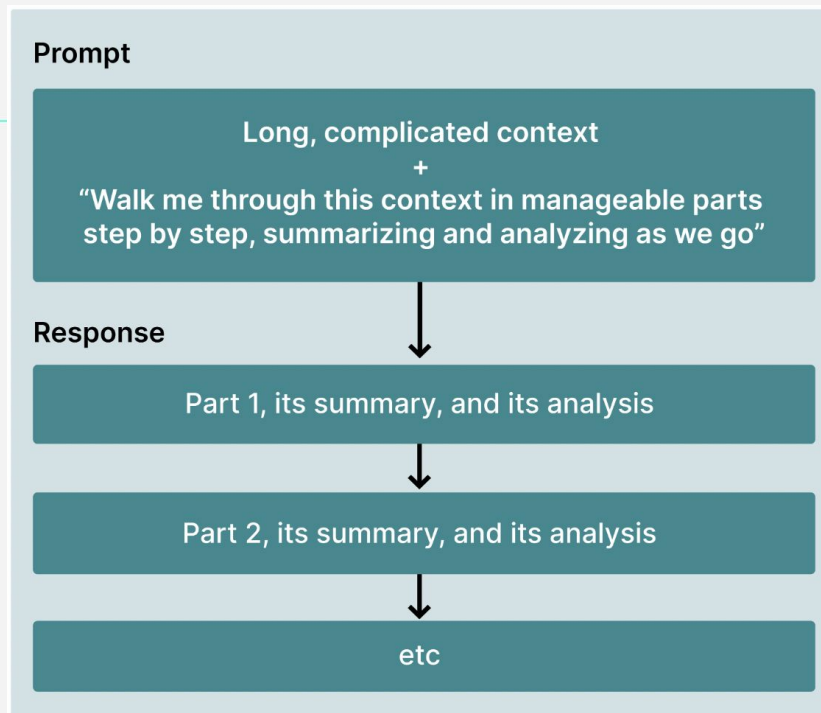
**Response**

Thought 1

Thought 2

Thought 3

Answer

# Thread-of-Thought

Virma has three bags, each of which fits five shirts. How many shirts can Virma fit in her bags?

Walk me through this context in manageable parts step by step, summarizing and analyzing as we go.
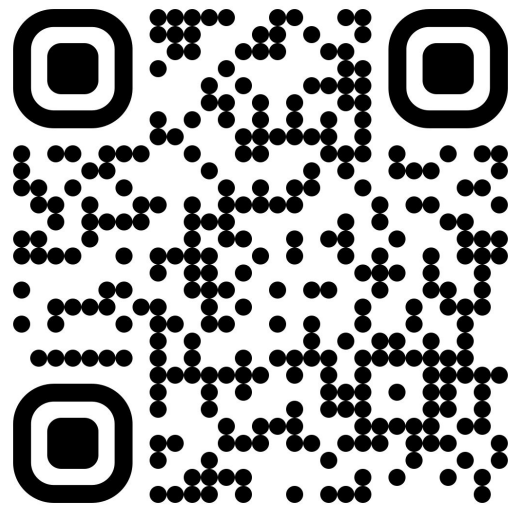
CIKLUM

**Prompt**

Long, complicated context
+
"Walk me through this context in manageable parts step by step, summarizing and analyzing as we go"

**Response**

Part 1, its summary, and its analysis

Part 2, its summary, and its analysis

etc

# CIKLUM

# Thank you!

# Share your feedback!

CIKLUM

Join our team

CIKLUM