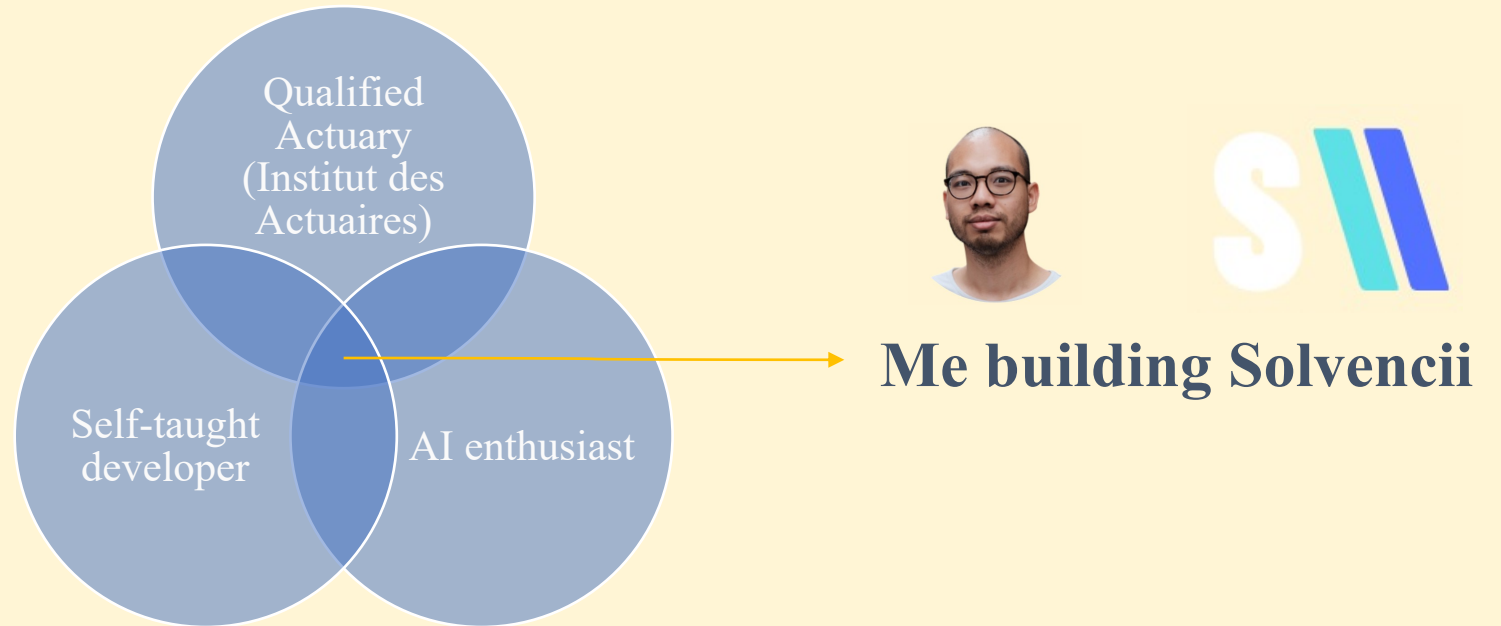Asian Actuarial Conference 2025

# Agentic AI in Asset and Liability Management (ALM)

**Duc Hien VU**

**SOLVENCII | Founder & CEO**

# About Me

Qualified Actuary (Institut des Actuaires)

Self-taught developer

AI enthusiast

**Me building Solvencii**

I'm not alone. I have 4 AI assistants working full-time @**Solvencii Copilot**.

Reach out to us directly at: **solvencii.fr** !!!
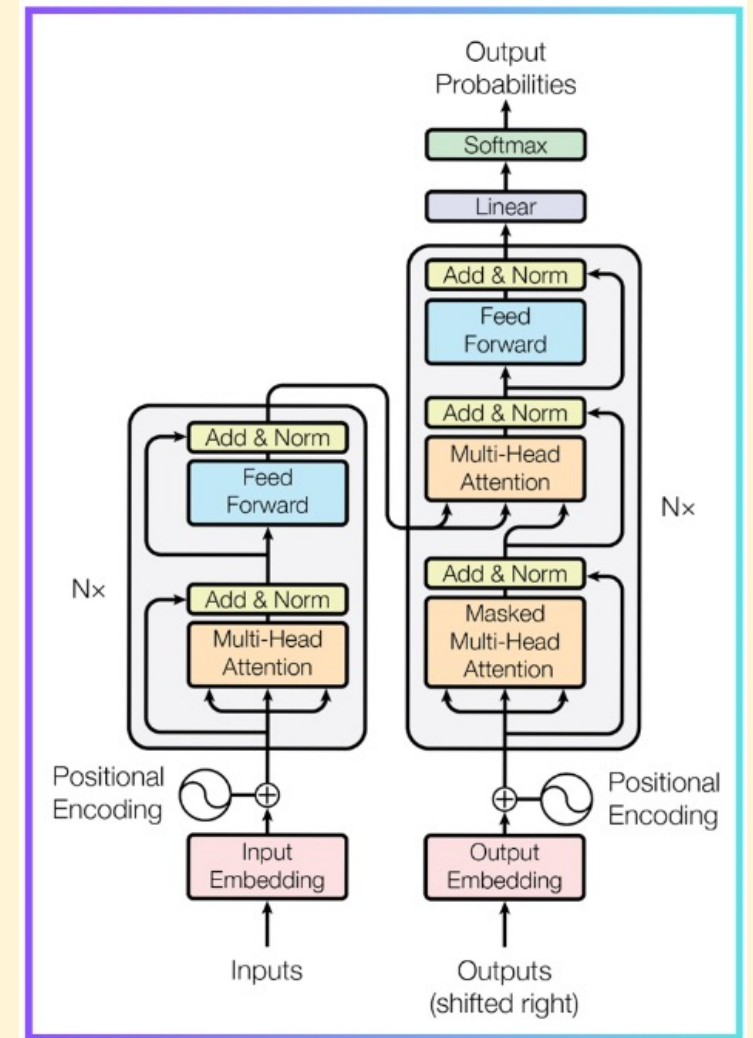
# Theoretical Foundations of LLMs

# Transformer Architecture

**A deep neural network framework that leverages self-attention mechanisms to process sequential data without relying on recurrent structures.**

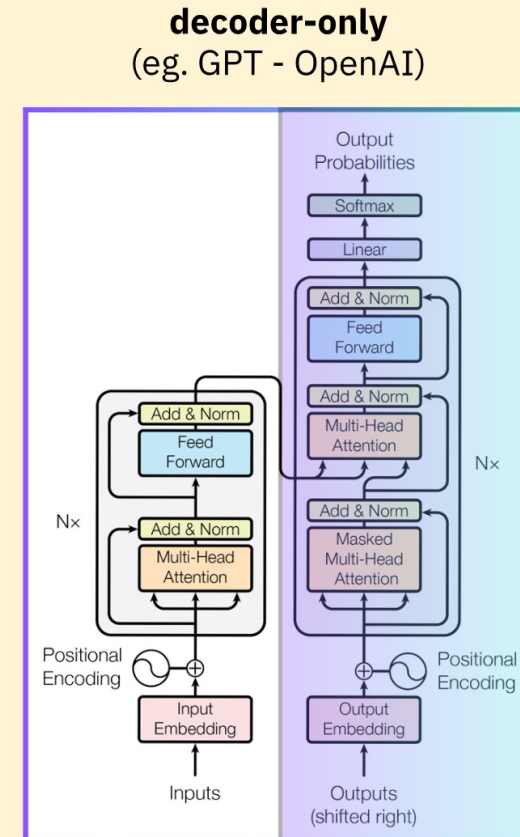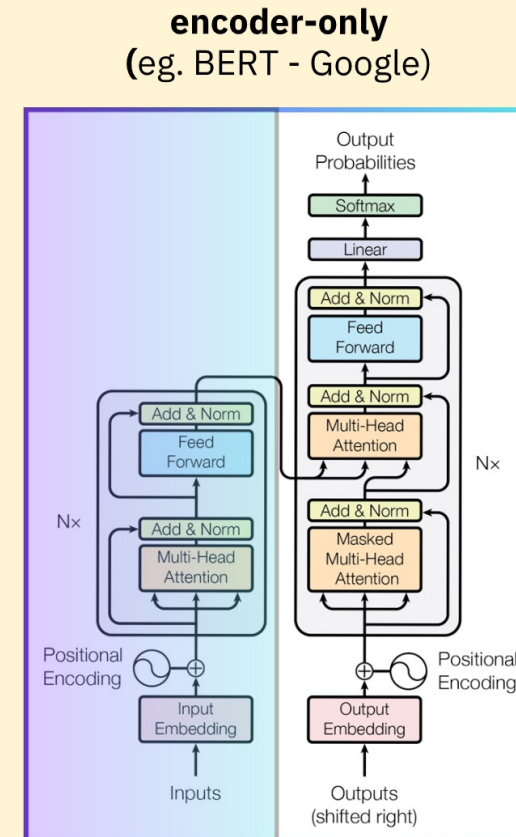It consists of stacked encoder and decoder layers, where:

- each encoder layer computes self-attention over input tokens to capture contextual relationships, and
- each decoder layer performs both self-attention and encoder-decoder attention to generate outputs.



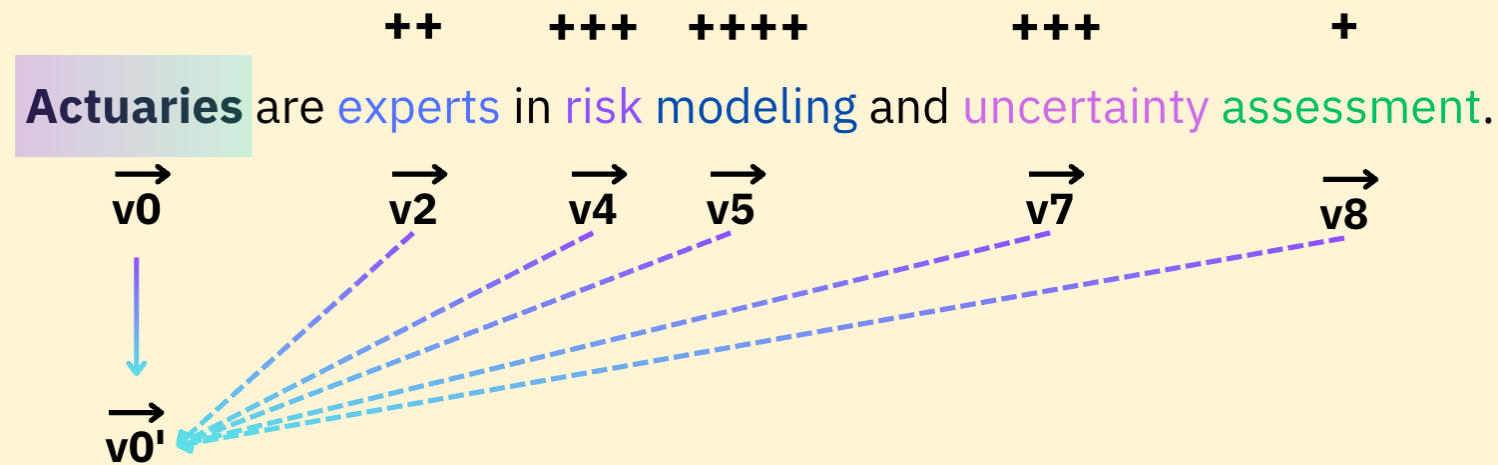*Attention is All You Need (Vaswani et al., 2017)*

# Variations of Transformer Architecture

- Encoder-only models focus on bidirectional context representation, ideal for understanding and classification tasks.

- Decoder-only models, are optimized for generative tasks by predicting subsequent tokens in a sequence, operating in a unidirectional manner.

- Encoder-decoder models combine both approaches to handle sequence-to-sequence tasks such as translation, summarization, and question answering

- Hybrid architectures that integrate elements of both self-attention and other neural structures have emerged.



**encoder-only**
(eg. BERT - Google)

**decoder-only**
(eg. GPT - OpenAI)

# Self-Attention mechanism

- Allows the model to capture the contextual meaning of the word **"Actuaries"** by attending to all other words in the sentence. First, the input words are converted into embeddings, which are projected into **query** $Q$, **key** $K$, and **value** $V$ vectors.



- For the word "Actuaries" the model computes a similarity score between its query vector $Q_{Actuaries}$ and the key vectors of all other words in the sentence using a dot product: $score(Q_{Actuaries}, K_i)$. These scores are passed through the **softmax** function to calculate attention weights (how much focus to give to each word). The final respresentation of "Actuaries" is computed as a weighted sum of the value vectors.

# Next Token Prediction

**Tokenization** breaks text into smaller, more manageable sub-word units

➔reduces the overall vocabulary size

➔allow the model to handle rare or out-of-vocabulary words more effectively

➔captures morphological nuances that full words might miss.

For a sequence of tokens $t_1, t_2, \ldots, t_n$, the **probability of a token $t_i$ given its preceding context** is modeled as:

$$P(t_i|t_{1:i-1}) = \frac{\exp(score(t_{1:i-1}, t_i))}{\sum_{t' \in V} \exp(score(t_{1:i-1}, t'))}$$

where $V$ is the vocabulary and $score(.)$ is a function learned by the network.

**Each generated token is appended to the existing sequence and fed back into the model as part of the input context to predict the next token.** This process is repeated iteratively until the model generates a stopping token or reaches the maximum token limit, completing the response.
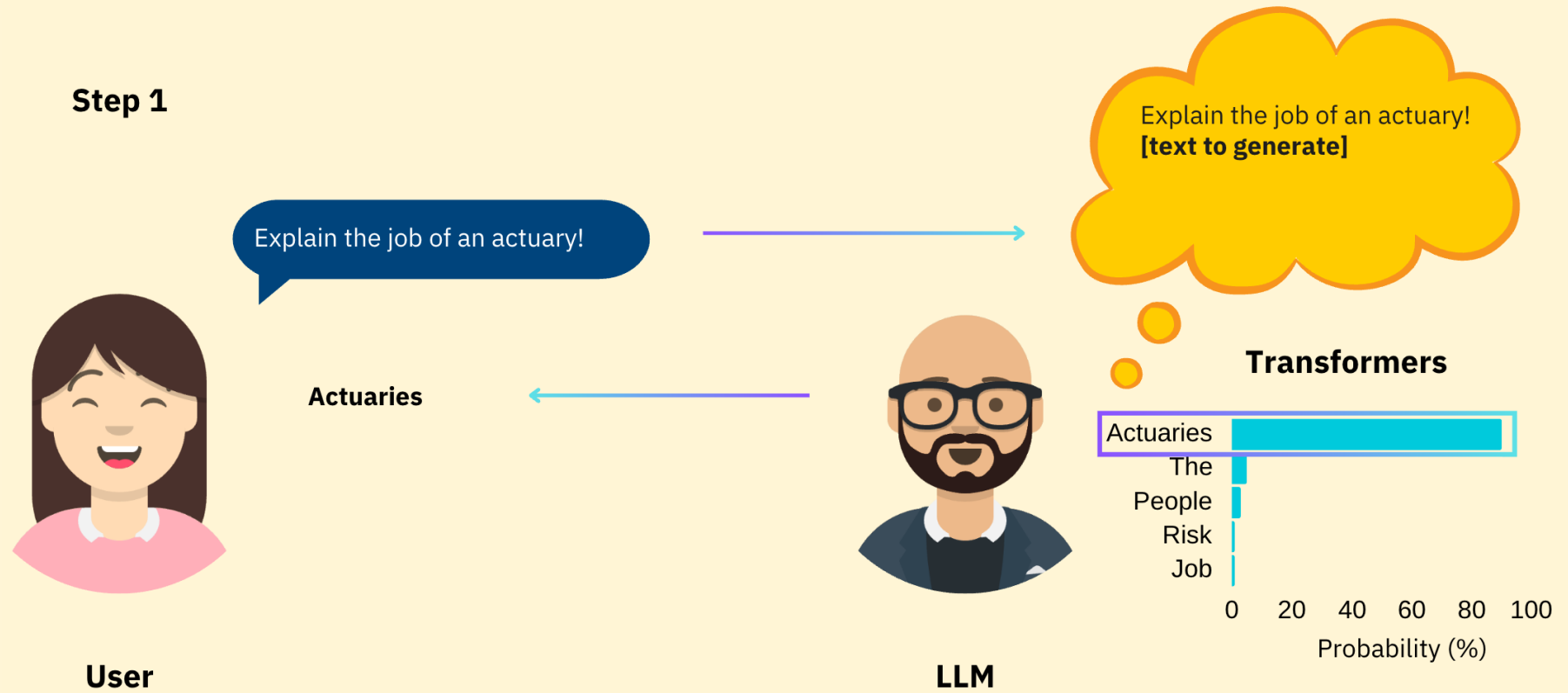
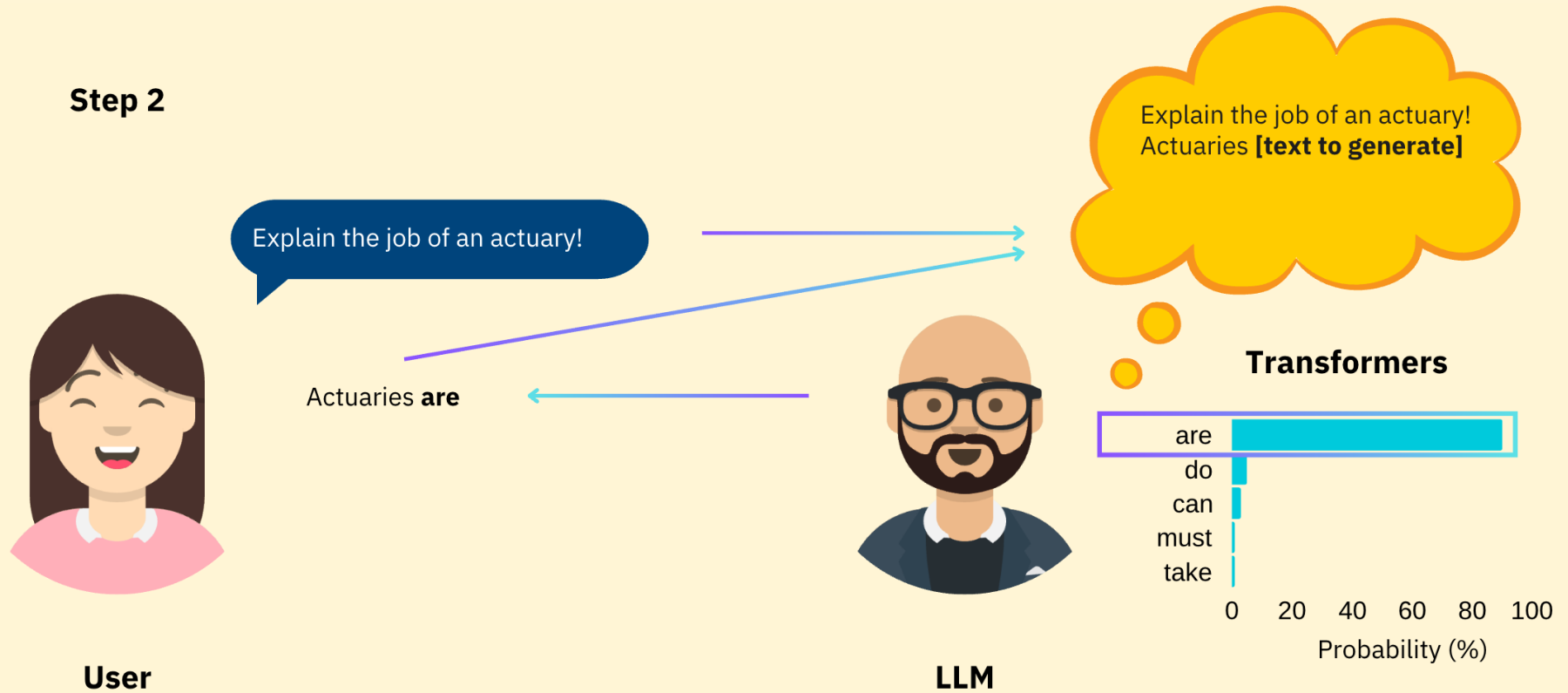# Next Token Prediction

**Step 0**

Explain the job of an actuary!
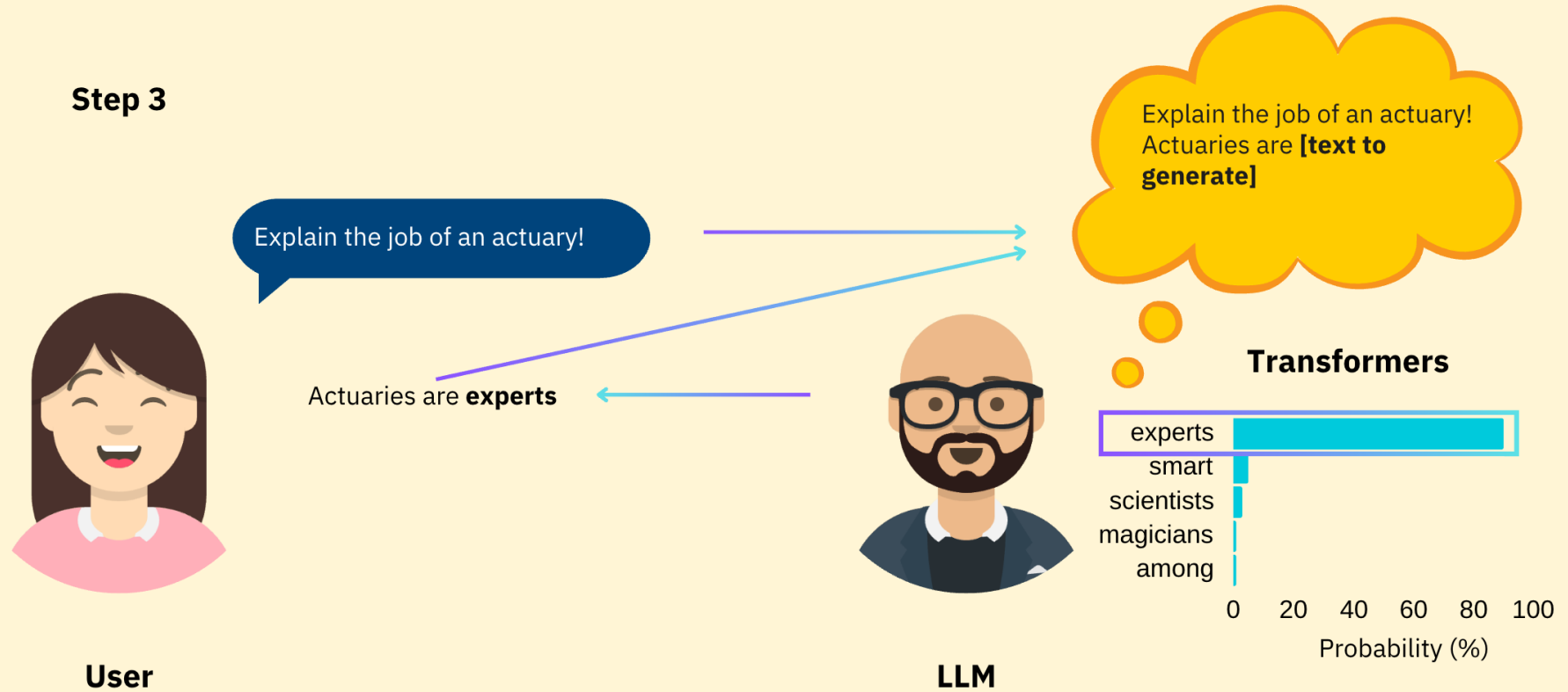
Explain the job of an actuary!
**[text to generate]**

**User**

**LLM**

# Next Token Prediction

# Next Token Prediction

# Next Token Prediction

**Step 3**

Explain the job of an actuary!

Explain the job of an actuary!
Actuaries are **[text to generate]**

Actuaries are **experts**

**Transformers**

experts
smart
scientists
magicians
among
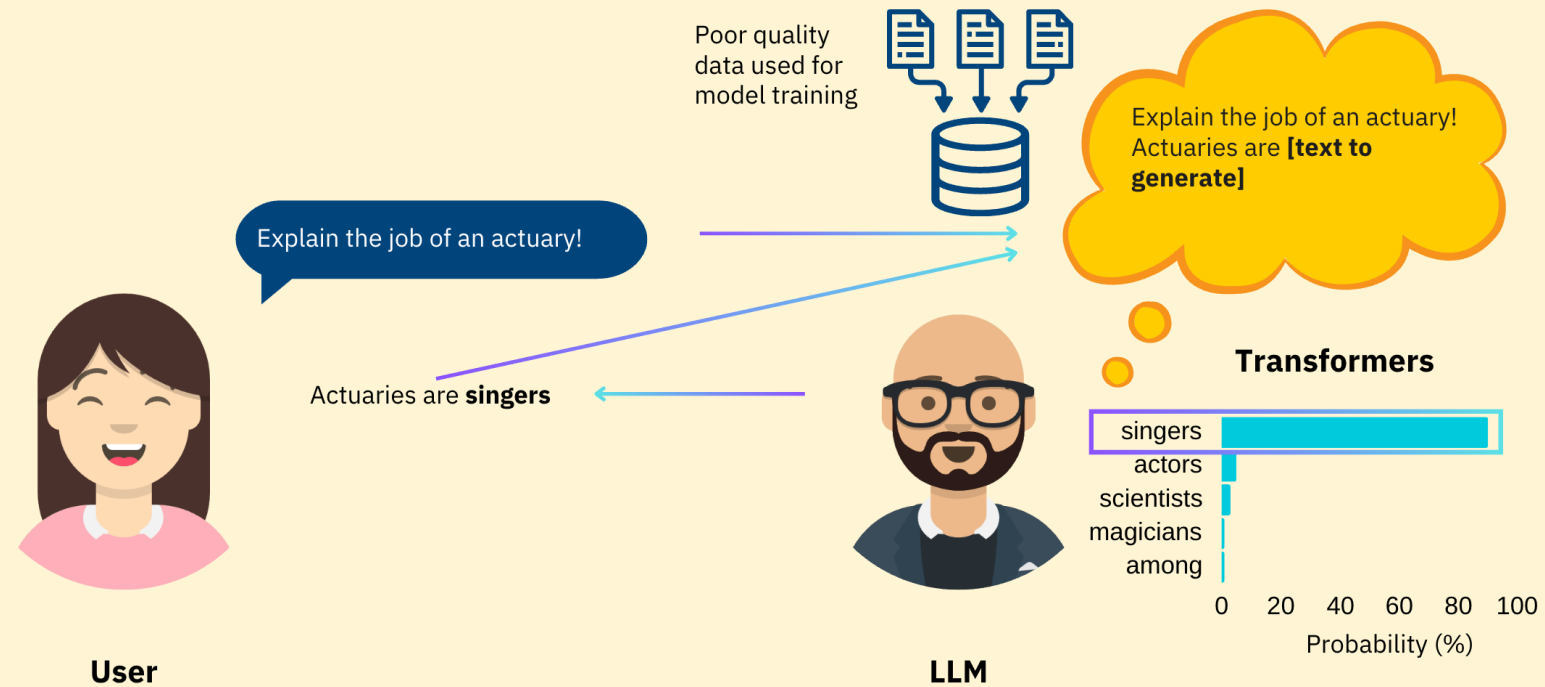
0  20  40  60  80  100
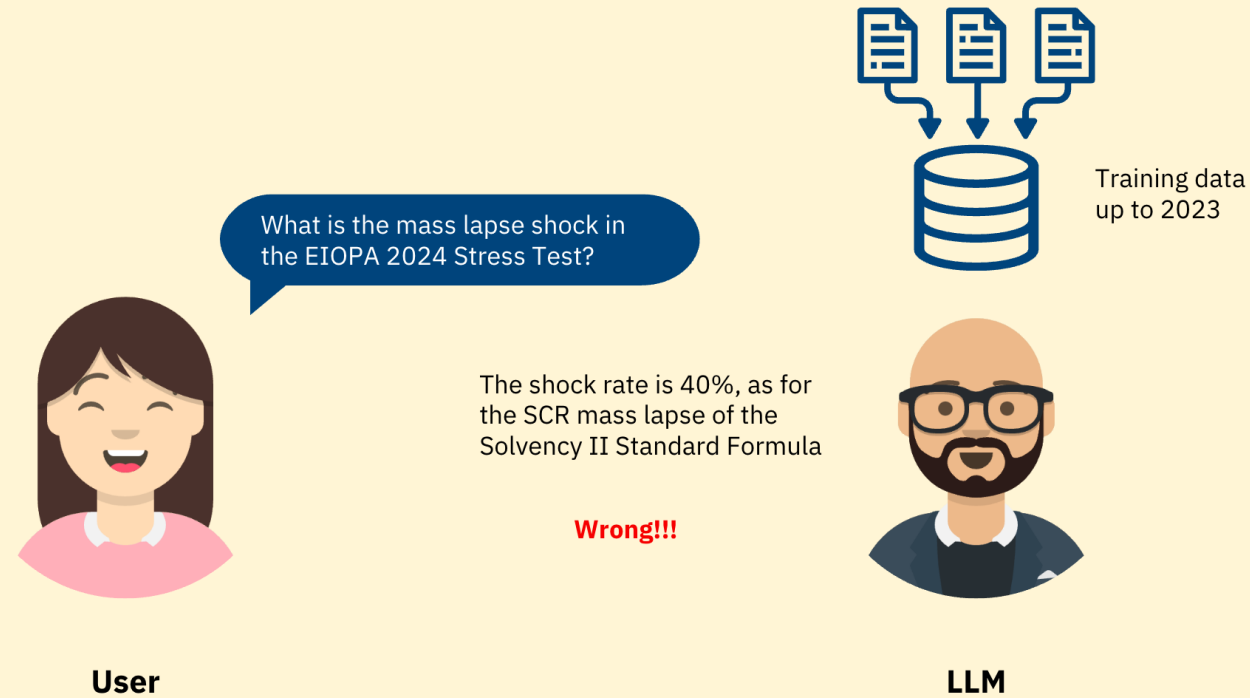Probability (%)

**User**

**LLM**

# Hallucination

- **Hallucination** in large language models (LLMs) refers to the phenomenon where the model generates factually incorrect, misleading, or nonsensical information that appears plausible.

- This occurs because **LLMs are trained to predict the likely next token based on patterns in the training data** rather than verifying the factual accuracy of the output.

- Hallucination is particularly problematic in professional fields like finance and actuarial science, where accuracy and reliability are crucial.

- The root causes of hallucination include **insufficient or biased training data, overgeneralization, lack of real-world grounding, and limitations in the model's ability to access or retrieve updated information**.

# Hallucination



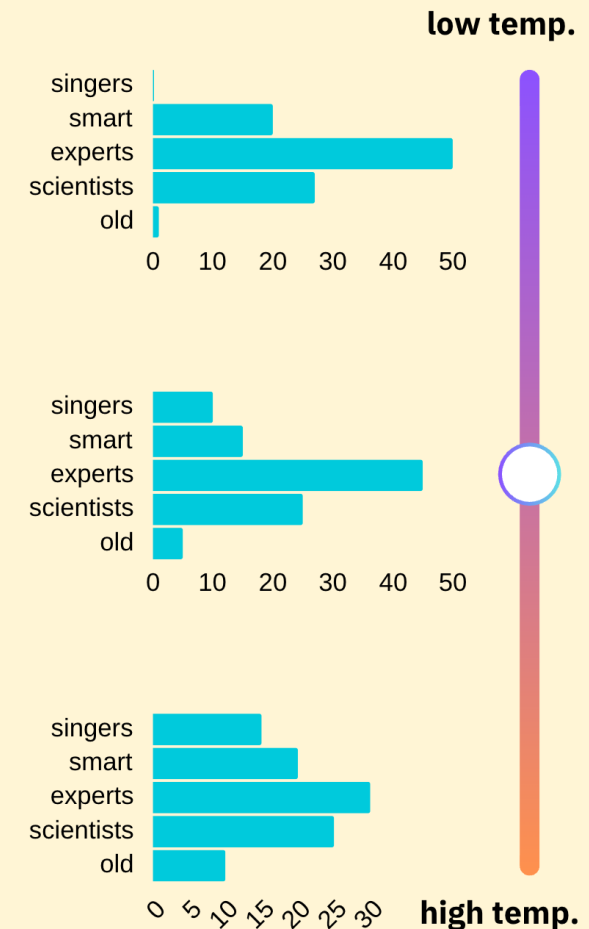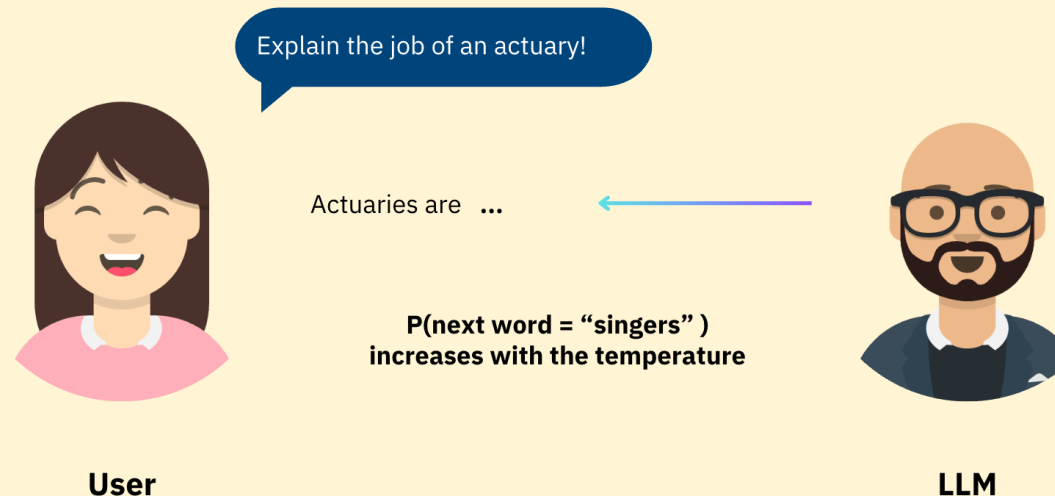*Hallucination due to poor quality data used for model training*

# Hallucination



*Hallucination due to limitations in the model's ability to access or retrieve updated information.*

# Temperature

Adjusting the temperature parameter can boost the quality of model outputs depending on the desired behavior.

Explain the job of an actuary!

Actuaries are ...

P(next word = "singers" )
increases with the temperature

**User**

**LLM**

**low temp.**

singers
smart
experts
scientists
old

0   10   20   30   40   50

singers
smart
experts
scientists
old

0   10   20   30   40   50

singers
smart
experts
scientists
old

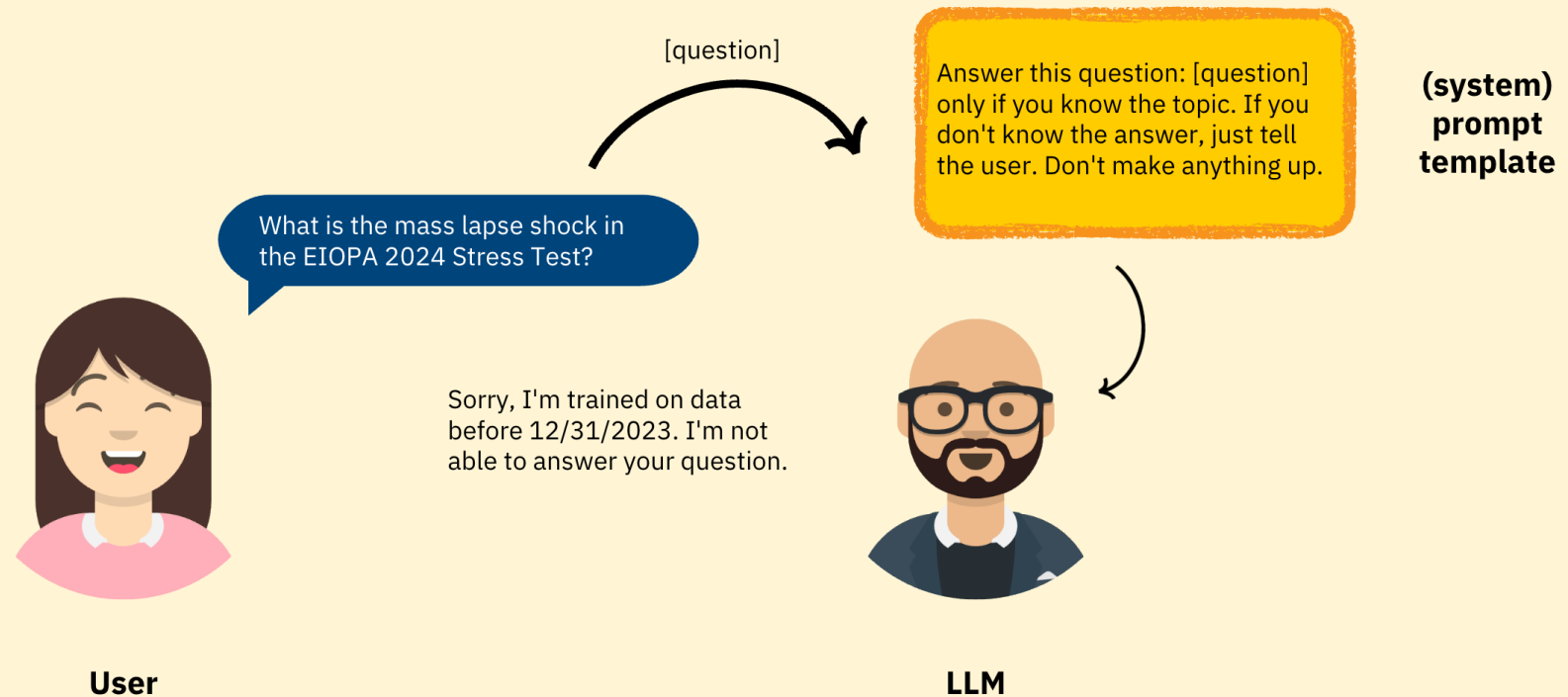0   5   10   15   20   25   30   **high temp.**

- A lower temperature sharpens the probability distribution, making the model more confident in its highest probability tokens and thus more deterministic in its output.
- A higher temperature flattens the distribution, leading to more diverse and random outputs.
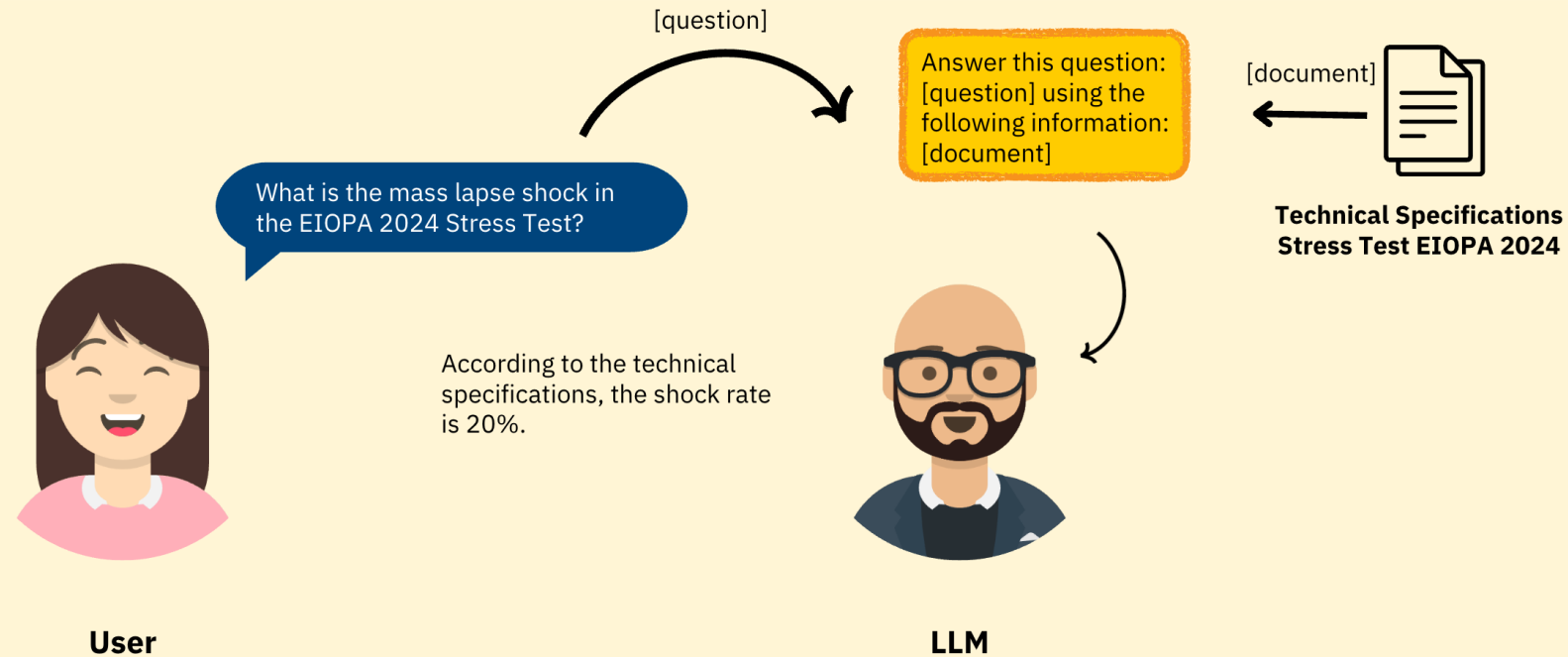
# Prompt Engineering

A common and effective way to reduce model hallucination is through **prompt engineering**. When a user submits a question, the input prompt is combined with a **system prompt template** that includes specific instructions designed to guide the model's behavior and improve response accuracy.

Returning to the example of hallucination caused by the model's limited access to updated information, one straightforward approach is to **explicitly instruct the model not to generate information it is unsure about**.

[question]

Answer this question: [question] only if you know the topic. If you don't know the answer, just tell the user. Don't make anything up.

**(system) prompt template**

What is the mass lapse shock in the EIOPA 2024 Stress Test?

Sorry, I'm trained on data before 12/31/2023. I'm not able to answer your question.

**User**

**LLM**

# Augmented Data & Context Window

If we still want the LLM to provide a meaningful answer despite its limited internal knowledge, one approach is to **inject augmented data into the system prompt** to give the LLM a more complete context.
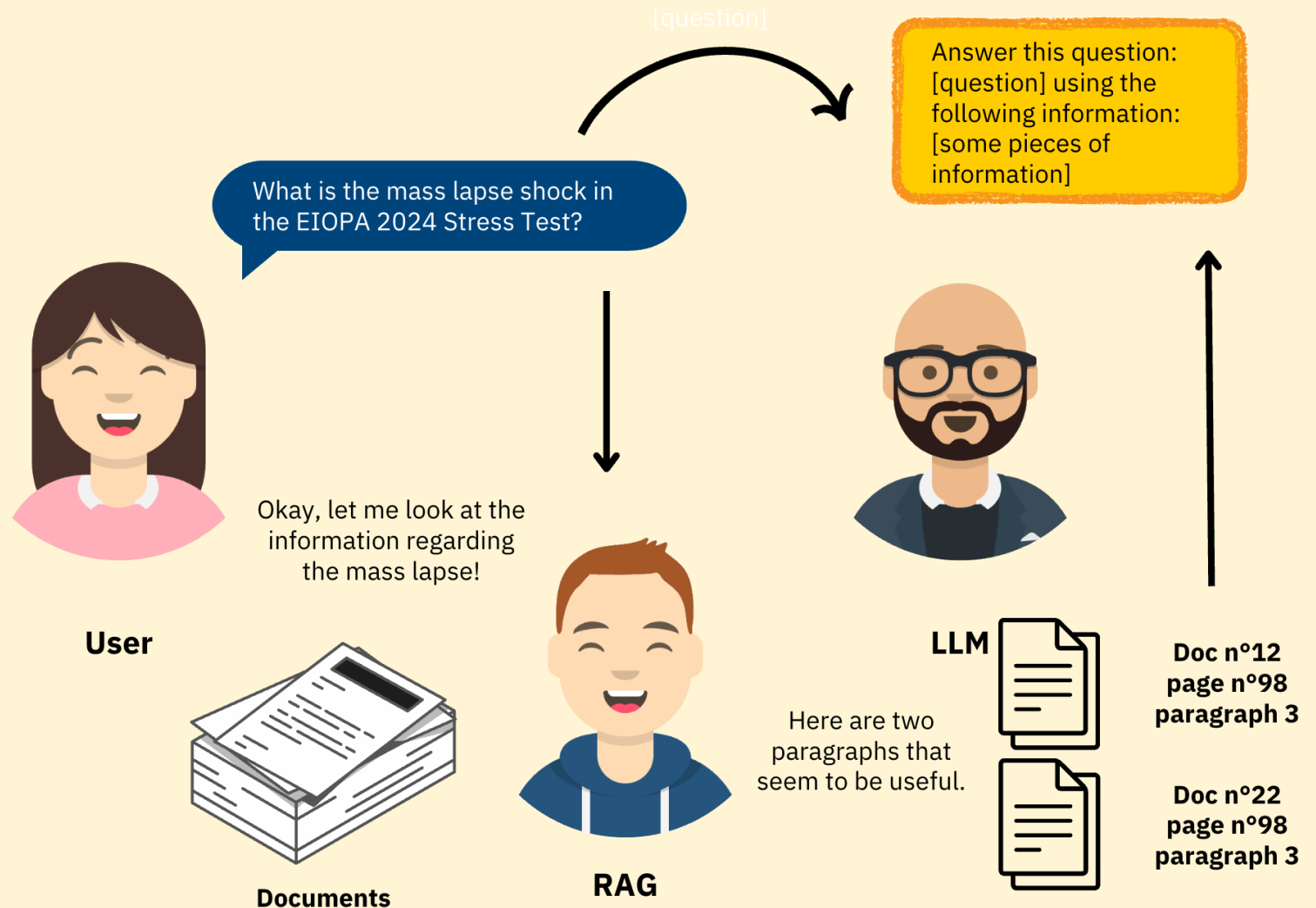


However, this is technically challenging because augmented data often resides in large documents, and the **context window** - which defines the maximum length of the prompt - is limited.

# Retrieval-Augmented Generation (RAG)

Even though context windows are expanding, filling the prompt with large amounts of raw data is often inefficient and costly.
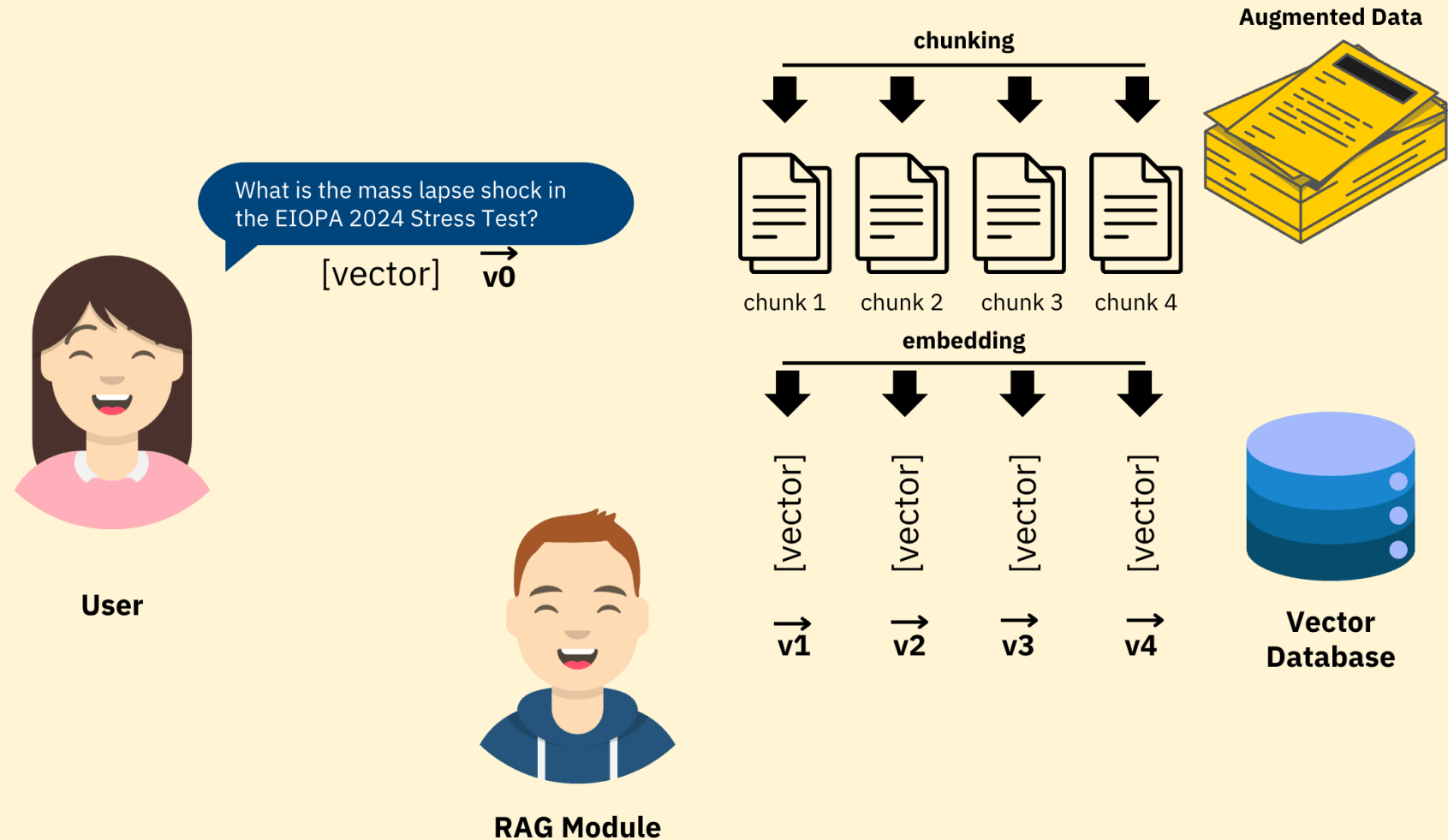
A smarter and more cost-effective solution is to **extract only the most relevant and useful pieces of information** and supply them to the model as context.

[question]

Answer this question: [question] using the following information: [some pieces of information]

What is the mass lapse shock in the EIOPA 2024 Stress Test?

Okay, let me look at the information regarding the mass lapse!

**User**

Here are two paragraphs that seem to be useful.

**LLM**

**Doc n°12 page n°98 paragraph 3**

**Doc n°22 page n°98 paragraph 3**

**Documents**

**RAG**

# Retrieval-Augmented Generation (RAG)

To enable efficient searching, text data is first converted into numerical representations called **embeddings**.

A vector database stores these embeddings in a structured format, allowing for efficient similarity-based searches.

What is the mass lapse shock in the EIOPA 2024 Stress Test?

[vector] $\vec{v0}$

**User**

**RAG Module**

chunking

chunk 1    chunk 2    chunk 3    chunk 4

embedding

[vector]    [vector]    [vector]    [vector]

$\vec{v1}$    $\vec{v2}$    $\vec{v3}$    $\vec{v4}$

**Augmented Data**

**Vector Database**

# Vector Search

When the user submits a query, it is embedded into a vector using the same embedding model.

The system performs a **vector search** to find the most similar vectors in the database using a distance metric like **cosine similarity** or **Euclidean distance**.

# Retrieval-Augmented Generation (RAG)

The closest matches (most semantically similar) are retrieved and passed to the LLM as part of the prompt context.

More advanced techniques are under development and testing to enhance RAG's performance and accuracy, such as **Graph RAG** and other hybrid retrieval approaches.

What is the mass lapse shock in the EIOPA 2024 Stress Test?

Answer this question: [question] using the following information: [chunk 1] and [chunk 4]

According to the technical specifications, the shock rate is 20%.

**User**

**LLM**

chunk 1

chunk 4

**RAG Module**

# Chain-of-Thought (CoT)

Example: **Prompt-based few-shot CoT (Inference-time technique)** → guide the model to think step-by-step by showing, in the system prompt, examples of step-by-step reasoning before the question.

By showing **structured examples**, we:
- Teach it **how to break down** the actuarial logic (mortality probabilities, discounting, summation).
- Implicitly define **the reasoning style**: short, precise, and actuarial.
- Encourage **step-by-step numeric reasoning**, reducing arithmetic or conceptual errors.

[question]

A 3-year term insurance of 100,000 is issued to a life aged 40.
$q_{40} = 0.002$, $q_{41} = 0.0025$, $q_{42} = 0.003$, and $i = 2\%$.
What is the EPV of benefits?

**User**

A: Let's think step by step.
1. Benefit amount = 100,000
2. Death can occur in year 1, 2, or 3.
3. Probabilities:
  - Year 1: $q_{40} = 0.002$
  - Year 2: $(1 - q_{40}) \times q_{41} = 0.998 \times 0.0025 = 0.002495$
  - Year 3: $(1 - q_{40})(1 - q_{41}) \times q_{42} = 0.998 \times 0.9975 \times 0.003 = 0.002985$
4. Discount factors:
  - $v^1 = 0.98039$, $v^2 = 0.96117$, $v^3 = 0.94232$
5. EPV = $100{,}000 \times [0.002 \times 0.98039 + 0.002495 \times 0.96117 + 0.002985 \times 0.94232]$
6. EPV = $100{,}000 \times 0.006346 = 634.6$
Answer: The EPV is approximately 634.6.

**LLM**

You are an actuary. Compute the expected present value (EPV) of benefits for each case.
Show your reasoning step by step before giving the final answer.

Example 1:
Q: A 1-year term insurance of 100,000 is issued to a life aged 40.
The annual mortality rate $q_{40} = 0.002$ and interest rate $i = 2\%$.
What is the EPV of benefits?

A: Let's think step by step.
1. Benefit amount = 100,000
2. Probability of death within the year = $q_{40} = 0.002$
3. Present value factor = $1 / (1 + i) = 1 / 1.02 = 0.98039$
4. EPV = $100{,}000 \times 0.002 \times 0.98039 = 196.08$
Answer: The EPV is approximately 196.08.

Example 2: ...

Now answer this question in the same reasoning style:
Q: [question]

# Tool Use

# Distillation

A model compression technique where a smaller, more efficient model (the *student*) is trained to replicate the performance of a larger, more complex model (the *teacher*).

Train the smaller model to match the teacher's output probabilities, hidden state representations, or attention patterns.
→ reduce the computational and memory requirements of LLMs while preserving most of the teacher model's accuracy and reasoning capabilities.

Months of training on massive datasets, using extensive compute resources on a lot of Nvidia H100 GPUs (cost $$$$$$$)

Direct teacher training (API calls) with much less computing resources (cost $)

Distillation

**The teacher**

**The student**

# AI Agents

System that can perceive its environment, process information, and take actions to achieve specific goals, often with a degree of autonomy and adaptability.

An AI agent typically consists of the following core components:

- **Perception** – The agent receives input from its environment, which could be user prompts, external data feeds, or system-generated signals.

- **Memory/State** – The agent maintains a state that tracks the current context and past interactions to inform future decisions.

- **Reasoning/Decision-Making** – The agent processes the input using LLMs or other models to evaluate possible actions and determine the best course of action.

- **Action/Execution** – The agent executes the chosen action, which could involve generating text, retrieving information, running calculations, or triggering an external process.

- **Feedback Loop** – The agent evaluates the outcome of its action and adjusts its future decisions based on the results, creating a learning loop.

# Building AI Agents

- **Define the Objective and Scope**
  - Identify the specific problem the agent should solve (e.g., interpreting ALM results, suggesting asset allocation strategies, etc.).
  - Set boundaries for the agent's behavior and capabilities (e.g., which data sources it can access, what decisions it can make).

- **Select the Core Model**
  - Choose an appropriate LLM (e.g., GPT, LLaMA) or fine-tune an existing model to suit the task.
  - Pre-train or fine-tune the model on domain-specific data (e.g., financial models, market data).

- **Build the Reasoning Engine**
  - Implement **chain-of-thought** reasoning to allow the agent to break down complex problems.
  - Use **tool-calling** and **function-calling** to allow the agent to interact with external systems (e.g., financial models, APIs).
  - Employ **planning** algorithms (e.g., tree search, reinforcement learning) to enable multi-step decision-making.

- **Integrate a Retrieval System**
  - Implement **RAG (Retrieval-Augmented Generation)** to allow the agent to access real-time market data, financial reports, and other external sources.
  - Fine-tune the retriever to optimize relevance and minimize noise.

- **Create a Feedback and Learning Loop**
  - Apply **Reinforcement Learning from Human Feedback (RLHF)** to refine outputs based on user evaluations.
  - Monitor the agent's performance and retrain periodically to adapt to changing market conditions or user needs.

- **Build an Interface**
  - Develop a user-friendly interface (e.g., chatbot or web app) to allow users to interact with the agent.
  - Enable multi-turn conversation handling to maintain context across user interactions.

# An effective AI Agent



**LLM reasoning**

**external tool interaction**

**adaptive learning**

a dynamic, goal-oriented system that can automate and optimize complex processes

# Asset and Liability Management

ALM modeling & ALM software

# Asset and Liability Management (ALM)

A core discipline within actuarial science that focuses on managing the interaction between an institution's assets and liabilities to optimize financial performance, minimize risk, and ensure long-term solvency.

# ALM Model Structure

Connects asset and liability models to simulate how assets and liabilities evolve together over time under different market scenarios.

- Applies dynamic asset allocation rules.
- Models reinvestment and disinvestment strategies.
- Simulates the impact
- of market shocks and
- economic cycles
- Models crediting and profit-sharing strategies

Models the behavior of financial markets, asset returns, interest rates, and credit spreads.
Uses stochastic models to simulate market dynamics.
Incorporates historical correlations between asset classes

Asset model

Interaction engine

Liability model

Projects future cash outflows based on policyholder behavior, mortality rates, claims experience, and regulatory requirements.

- Uses actuarial models (e.g., survival models, annuity factors) to simulate liability cash flows.
- Accounts for optionality in insurance contracts (e.g., policyholder lapses, surrender options).

# Challenges in traditional ALM software

- **Data Preparation and Input Complexity:** ALM models require the ingestion of diverse datasets—from market indices and interest rates to demographic statistics and policy details. Standardizing these data sources into a model-ready format is nontrivial and often requires manual intervention, increasing the potential for errors.

- **Algorithmic Sophistication in ALM:** ALM algorithms involve complex processes such as asset allocation, investment strategies, crediting methodologies, and profit-sharing mechanisms. Each component demands precise calibration and robust mathematical modeling, which can overwhelm traditional computational systems, especially during extensive scenario simulations.

- **Output Interpretation and Decision-Making:** The output of traditional ALM models is often a vast array of scenarios and risk assessments. Interpreting these outputs requires deep domain expertise and can be time-consuming, reducing the effectiveness of decision-making under tight time constraints.

# Integrating AI Agents into ALM systems

**User Experience:**

- act as intelligent assistants, helping users navigate the different components of ALM software and interpret documentation.

- answer user queries, explain model outputs, and suggest improvements by leveraging domain knowledge from both the internal ALM model and external financial and regulatory sources available on the internet.

# Integrating AI Agents into ALM systems

**Automating Data Ingestion:** automate the process of collecting, cleaning, and standardizing data from diverse sources, including financial market data, actuarial reports, economic forecasts, and internal balance sheets, and updating model inputs in real time. → Reduce manual errors and ensures that the data is structured in a format compatible with ALM models ; Enable real-time adjustments to ALM strategies, improving the responsiveness of the model to changing market conditions.

# Integrating AI Agents into ALM systems

**Algorithmic Enhancement:** assist in fine-tuning ALM model parameters by employing optimization algorithms and machine learning techniques. (eg. reinforcement learning can be used to identify optimal asset allocation, investment strategies, systematic gain & loss realization, …)

# Integrating AI Agents into ALM systems

**Output Synthesis:** distill insights from complex simulation outputs.

- Identify patterns in large datasets,
- Generate natural language summaries of scenario outcomes, and provide actionable recommendations.

# Case Study

Solvencii Copilot in Solvencii Lab

# Solvencii Lab & Solvencii Copilot

**Solvencii Lab** is a SaaS platform for ALM modeling that integrates multiple modules for **data processing**, **simulation**, and **analysis** within a unified environment. Built with modern web technologies, Solvencii Lab offers a flexible and scalable architecture, making it easier to leverage state-of-the-art AI models and cutting-edge software techniques compared to traditional ALM software.

**Solvencii Copilot** is an advanced AI assistant embedded within Solvencii Lab, designed to extend the capabilities of traditional ALM software by integrating deep financial and actuarial expertise with state-of-the-art language models and tool orchestration. More than just a chatbot, Solvencii Copilot functions as a strategic co-pilot for actuaries, providing insights, automation, and enhanced decision-making support.

Both **Solvencii Lab** and **Solvencii Copilot** are openly and freely accessible at **solvencii.fr**. Users can explore the platform's full range of ALM modeling functionalities and experience the capabilities of Solvencii Copilot directly.

# Key capabilities of Solvencii Copilot



**Broad expertise in quantitative finance, actuarial science, and risk management**



**Retrieval-Augmented Generation (RAG) using external PDF files, with a fact-checking feature**



**Recommending Solvencii Lab's components and functionalities tailored to user queries**



**Conversational memory, self-reflection, and self-reasoning for enhanced problem-solving**

# Illustration

Solvencii Copilot is still evolving alongside the rapid advancements in LLMs and AI agent frameworks.

Its capabilities are expected to grow as AI models become more powerful, context windows expand, and multimodal capabilities improve.

Example scenario illustrating how Solvencii Copilot could assist a user within Solvencii Lab:

I wanna know the major impacts of a rise in interest rates on the profitability and the solvency of a life insurer.

(Synthetic answer based on underlying LLM data and Solvencii's documentation, fetched using RAG techniques): [•••]

Can you perform a modeling test on a fictional life insurer representing the average of the French life insurance market, to support your answer?

(Automatically browses the internet and downloads SFCR reports of the 10 biggest life insurers in France by market share in 2025 to construct the fictional insurance company): I've constructed the fictional insurer. Please specify the interest rate shock parameters.

I'd like to apply a shock of +200 basis points to the baseline yield curve provided by EIOPA for the EU in 2025.

Prepares a stochastic run for the test and presents a ready-to-use simulation): Here's the stochastic simulation setup. You can adjust parameters using Solvencii Lab UI if necessary.

That's good. Please launch the stochastic simulation, analyze the model outputs, and create a PowerPoint summarizing the results.

# Multi-agent architecture

- **Model diversification:** By leveraging different LLM models available on the market, Solvencii Copilot can match each task with the most suitable model. Even if we stick with a single model, we can assign different roles to various AI assistants using prompt templates.

- **Dynamic Orchestration:** Solvencii Copilot can dynamically decide which model to activate for a given task, optimizing for both **performance** and **cost**. For complex calculations or scenario analysis, a high-performing model may be required. For simpler tasks like data formatting or interpretation, a smaller, more efficient model may be sufficient.

- **Collaborative Problem Solving:** In some cases, instead of relying on a single large model, Solvencii Copilot can enable **collaboration between mid-sized models**. This mirrors the real-life scenario where a team of actuaries might work together—each contributing specialized knowledge to solve a complex ALM task.

- **Redundancy and Robustness:** A multi-agent setup also enhances system reliability. If one model underperforms or fails to deliver a satisfactory result, another model can take over or refine the output, ensuring consistent and accurate responses

# Behind the scenes
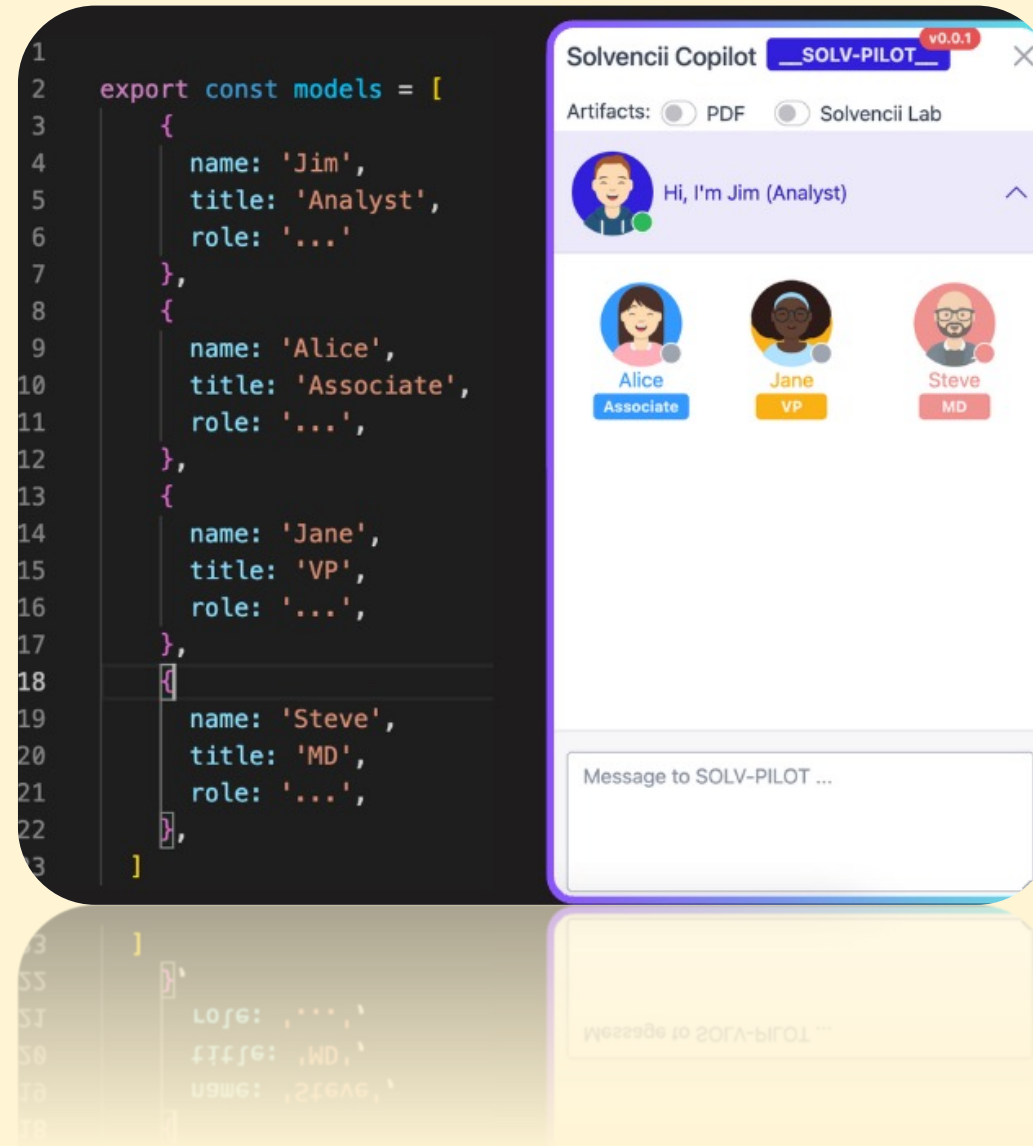
The making of Solvencii Copilot

# Model & Framework

# Multi-agent architecture

# Prompt Template, History & Memory

**prompt**

```
24  export const getAnswer = async (model, question, history) => {
25    let answer = ''
26    const systemTemplate =
27      `You are ` +
28      model.name +
29      `, ` +
30      model.title +
31      ` at Solvencii Copilot team (SOLV-PILOT), an AI assistant for
32      question-answering tasks on financial, actuarial and risk management
33      topics. You can interact with Solvencii Lab, a software for Asset and
34      Liability Mangement. Use your knowledge and this conversation history: `
35      history +
36      ` ` +
37      model.role

39    const prompt = ChatPromptTemplate.fromMessages([
40      ['system', systemTemplate],
41      ['human', '{input}'],
42    ])
43    const chain = prompt.pipe(llm).pipe(parser)
44    try {
45      answer = await chain.invoke({ input: question })
46    } catch (e) {
47      console.log(e)
48    }
49    return answer
50  }
```

```
questions = [
    'question n°1',
    'question n°2',
    'question n°3'
]

answers = [
    {name: 'Jim', message: 'réponse n°1'},
    {name: 'Alice', message: 'réponse n°2'},
    {name: 'Jane', message: 'réponse n°3'},
]

const history = answers.map(
    (a, i) =>
    'user: ' +
    questions[i] +
    '. ' +
    a.name +
    ': ' +
    a.message
)
```

```
output = await getAnswer(model, prompt, history.slice(-10).join(' '))
```

```
setAnswers((a) => [...a, { name: copilot.name, message: output }])
```

45

# RAG

# Thank you!

Do you have any questions ?

Not a chatbot this time

duc-hien.vu@solvencii.com
contact@solvencii.com

**Solvencii.com | Solvencii.fr**

https://www.linkedin.com/company/solvencii/