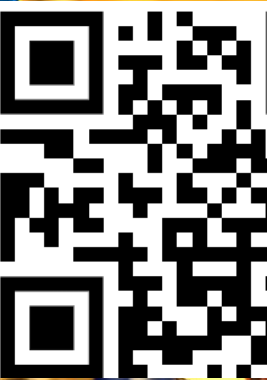


GenAI unpacked: Beyond Basics

<https://www.linkedin.com/in/damirdobric/>

Dr. Damir Dobric

Lead Software Architect daenet GmbH / ACP Digital
Microsoft Regional Director,
Most Valuable Professional: AI



AGENDA

- Tokens
- Embeddings
- Vector DBs
- RAG
- Function Calling
- Semantic Kernel



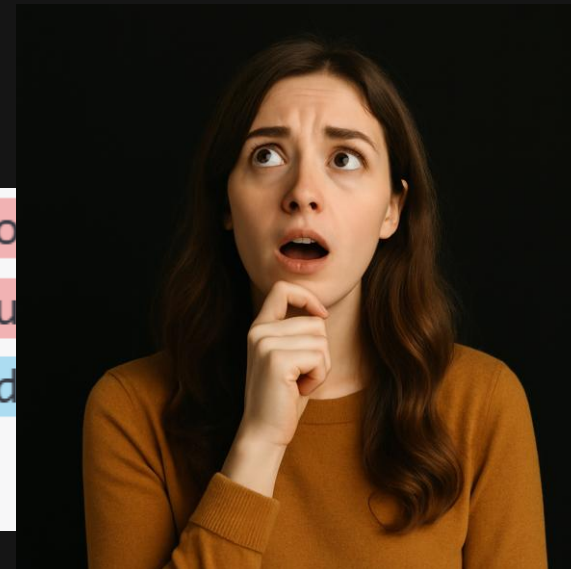
TOKENS



What are Tokens?

- 1 token \sim 4 chars in English
- 1 token \sim $\frac{3}{4}$ words
- 100 tokens \sim 75 words
- Byte Pair Encoding (Gage,1994): [Wikipedia](#)
- What are tokens and how to count them?
- Token Pricing: [Pricing \(openai.com\)](#)

OpenAI's large language models process text using tokens, which are continuous sequences of characters found in a set of text. The models learn to understand the statistical relationships between these tokens, and excel at producing the next token in a sequence of tokens.



DEMO

- Tokenizer
- Creating tokens in C#

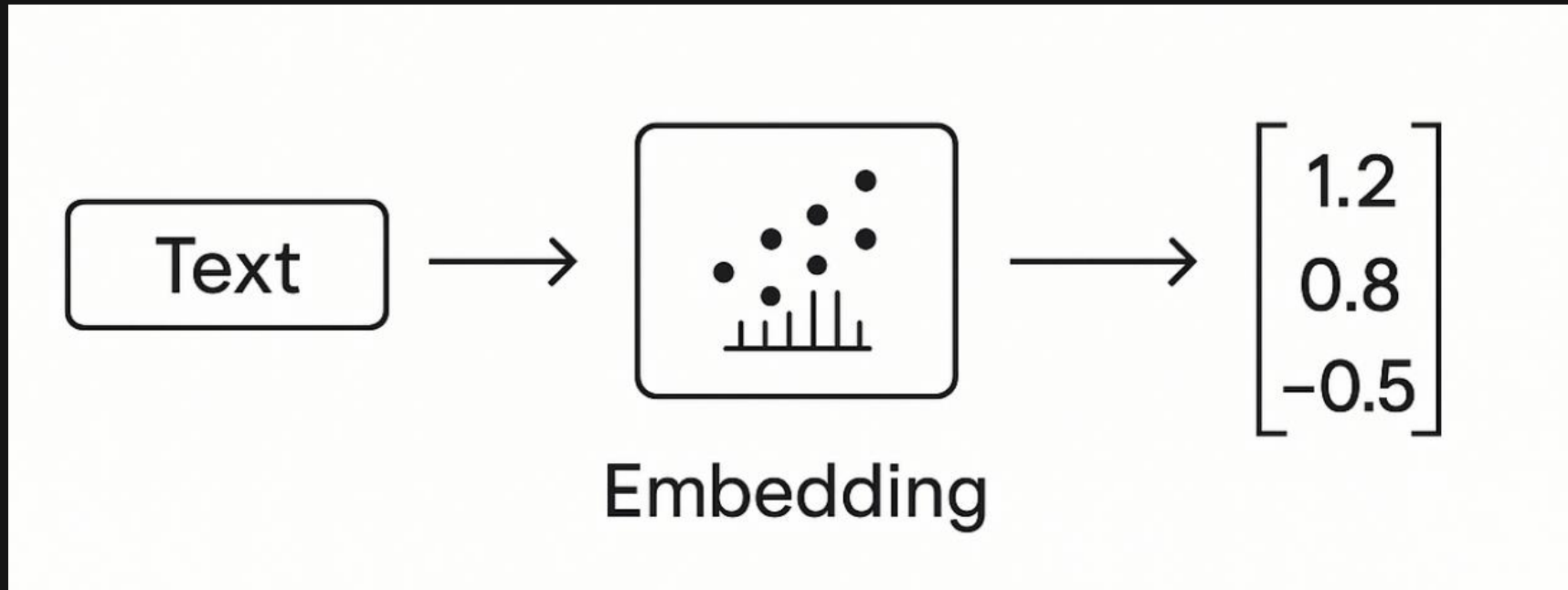


EMBEDDINGS

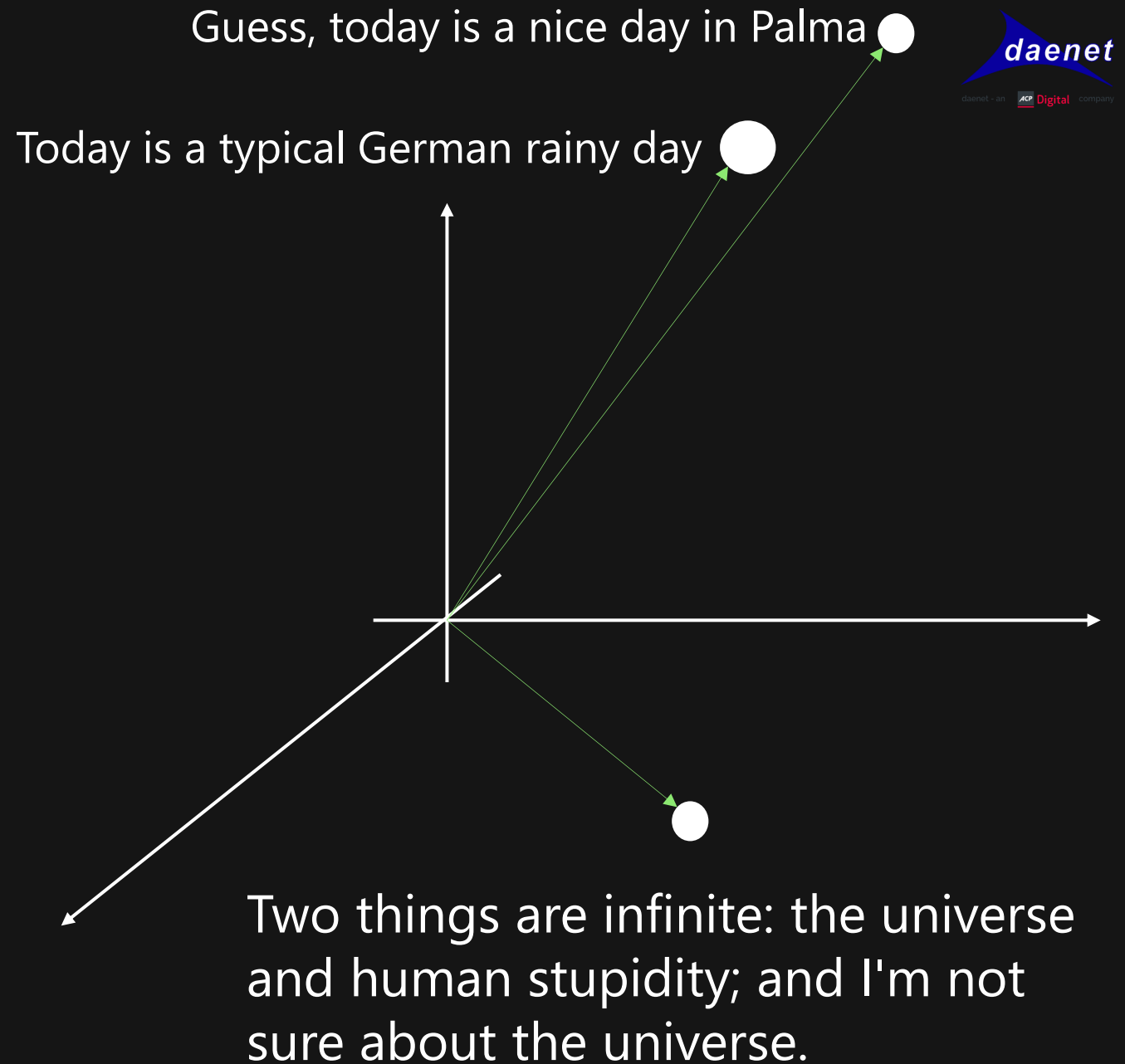


Embeddings

- Making a vector from text



Embedding Models



Similarity Between Multidimensional Vectors

- Dot Product
- The Norm
- Cosine Similarity

$$A \cdot B = a_1 \cdot b_1 + a_2 \cdot b_2 + \dots + a_n \cdot b_n$$

$$\|\mathbf{A}\| = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2}$$

$$\mathbf{A} \cdot \mathbf{B} = \|\mathbf{A}\| \|\mathbf{B}\| \cos(\theta)$$

When to use Embeddings?

- Semantic Search
- Classification
- Clustering
- Recommendations
- . . .

DEMO

- Consuming Embedding Model with HTTP/Post
- Creating Embeddings in C#



VECTOR DATABASES



Vector DBs

- SqlServer 2025 Native Vector Search
- Quadrant
- CosmosDB
- ...

SqlServer 2025 Native Vector Search

- SqlServer 2025 Native Vector Search
 - In Azure
 - On-Prem
- [SQL Server Native Vector Search for .NET Developers](#)

```
CREATE TABLE test.Vectors
(
    [Id] INT IDENTITY(1,1) NOT NULL,
    [Text] NVARCHAR(MAX) NULL,
    [VectorShort] VECTOR(3) NULL,
    [Vector] VECTOR(1536) NULL
);
```


DEMO

- SQL Server Native Vector Search

SQL Server Native Vector

```
DECLARE @v1 VECTOR(2) = '[1,1]';  
DECLARE @v2 VECTOR(2) = '[-1,-1]';  
  
SELECT  
    VECTOR_DISTANCE('euclidean', @v1, @v2) AS euclidean,  
    VECTOR_DISTANCE('cosine', @v1, @v2) AS cosine,  
    VECTOR_DISTANCE('dot', @v1, @v2) AS negative_dot_product;
```



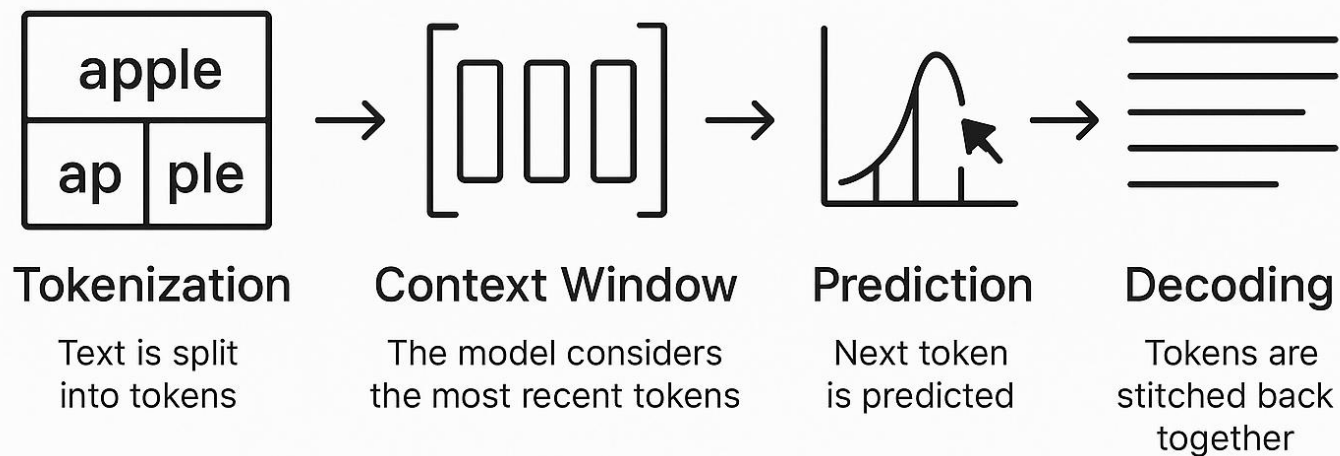
COMPLETION MODELS

How does all this work?



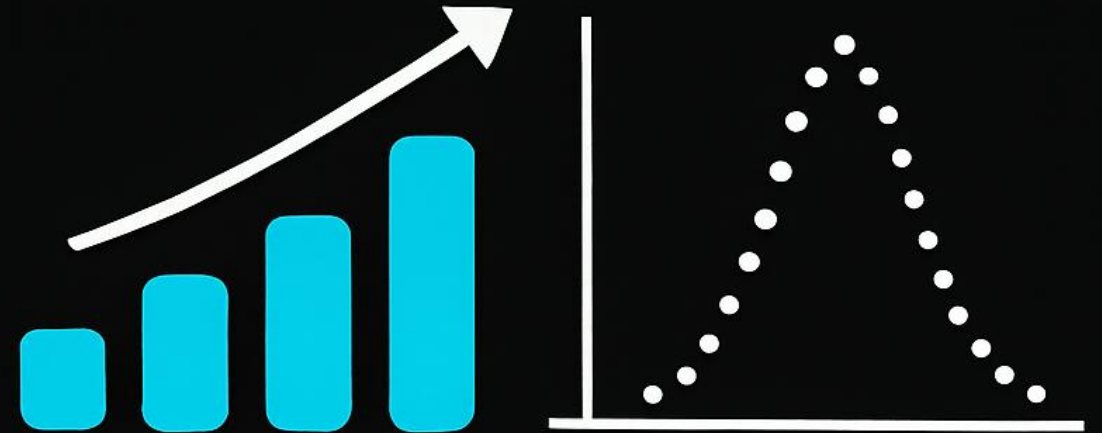
Completion Models

Text Completion Flow



DEMO

- C# Application
- Using OpenAI nuget package
- Executing ChatCompletions
- Presenting Probabilities

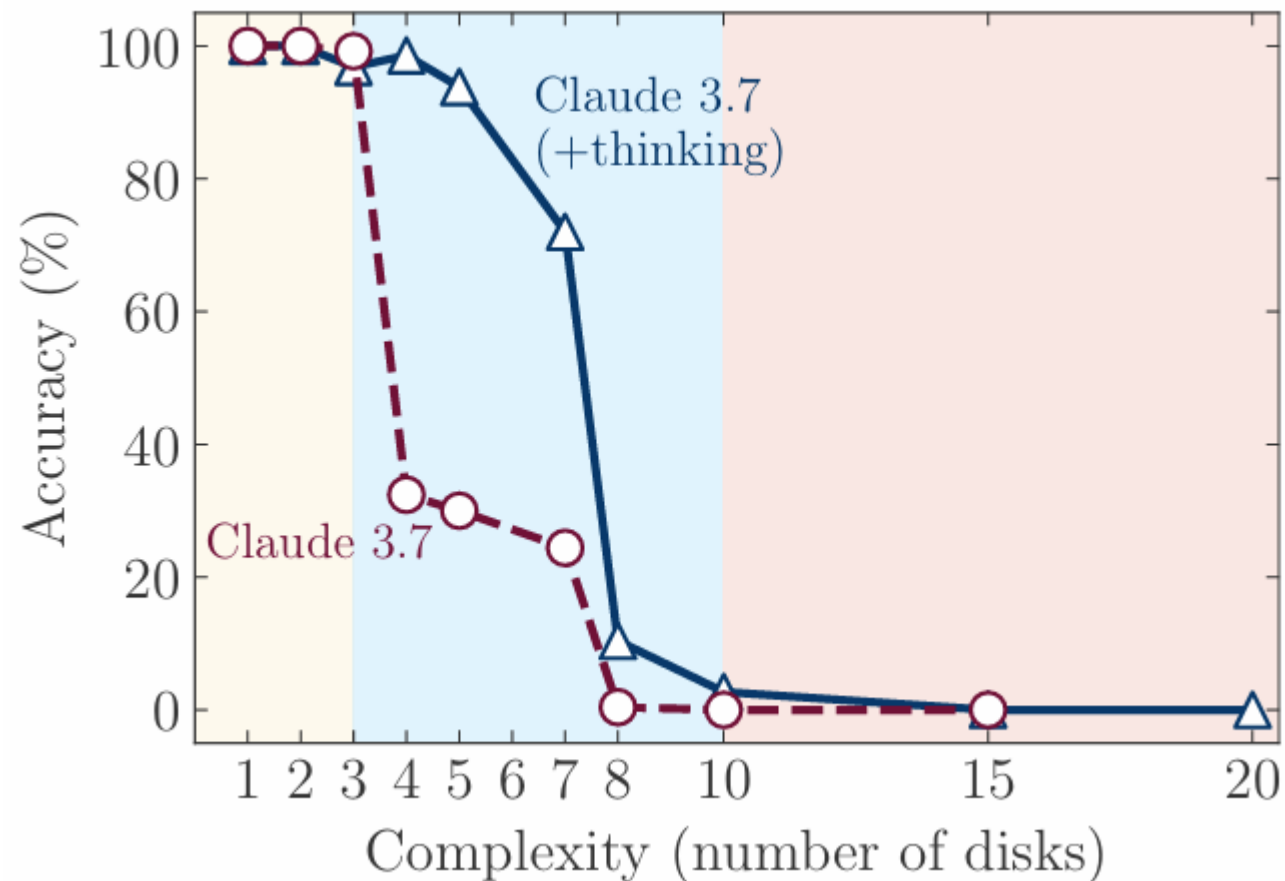


A Scientific look under the hub



Model Performance and Illusion of Tinking

[LiveBench](#)



ml-site.cdn-apple.com/papers/the-illusion-of-thinking.pdf

RAG

Retrieval Augmented Generation

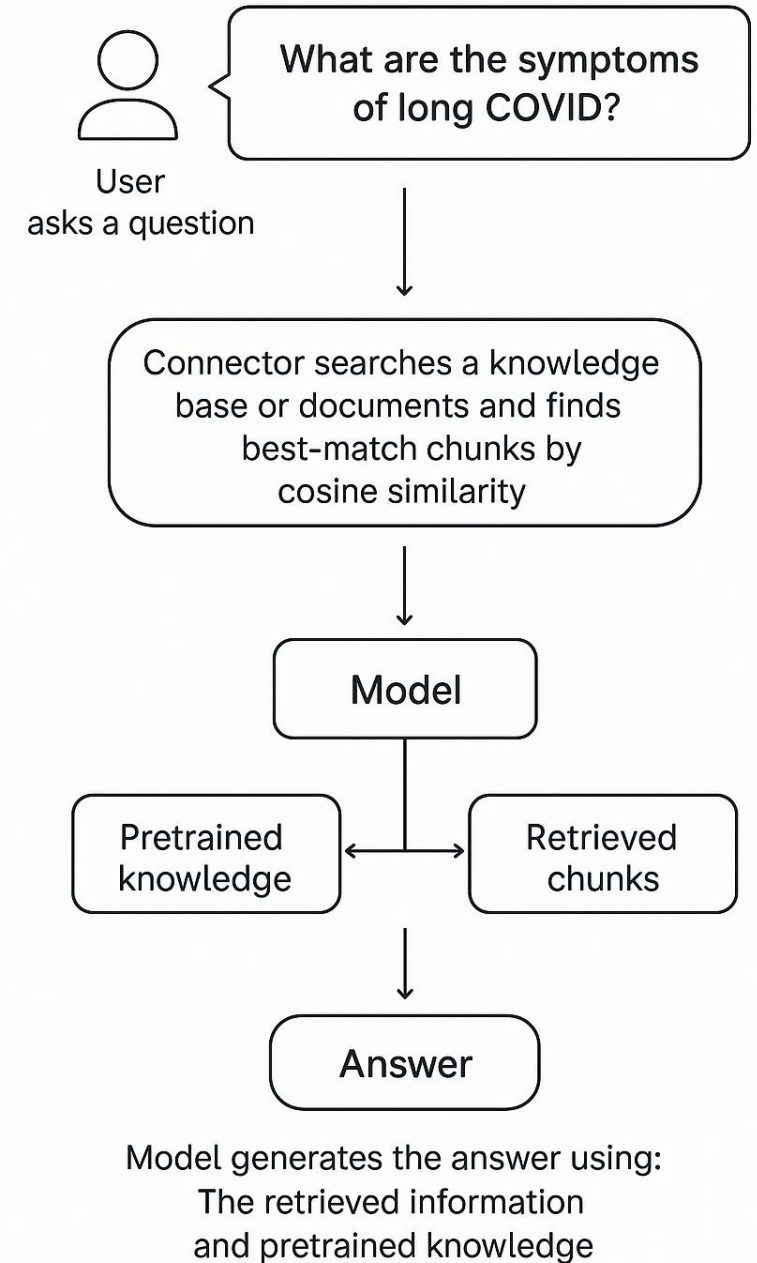


RAG

Retrieval-Augmented Generation

- Extending the model knowledge
- Combines information retrieval with text generation

<https://arxiv.org/abs/2005.11401>



DEMO (c#)

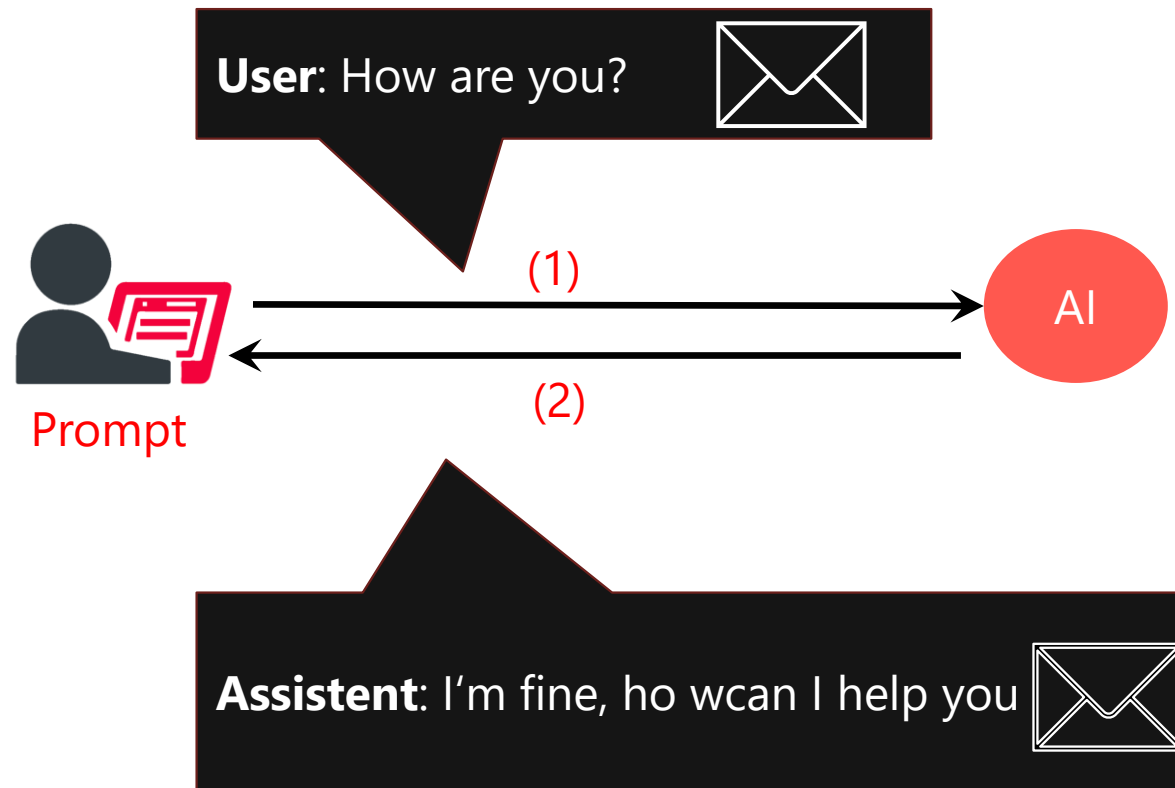
- Creating chunks
- Creating embeddings
- InMemory Vector DB
- Calculating Similarity



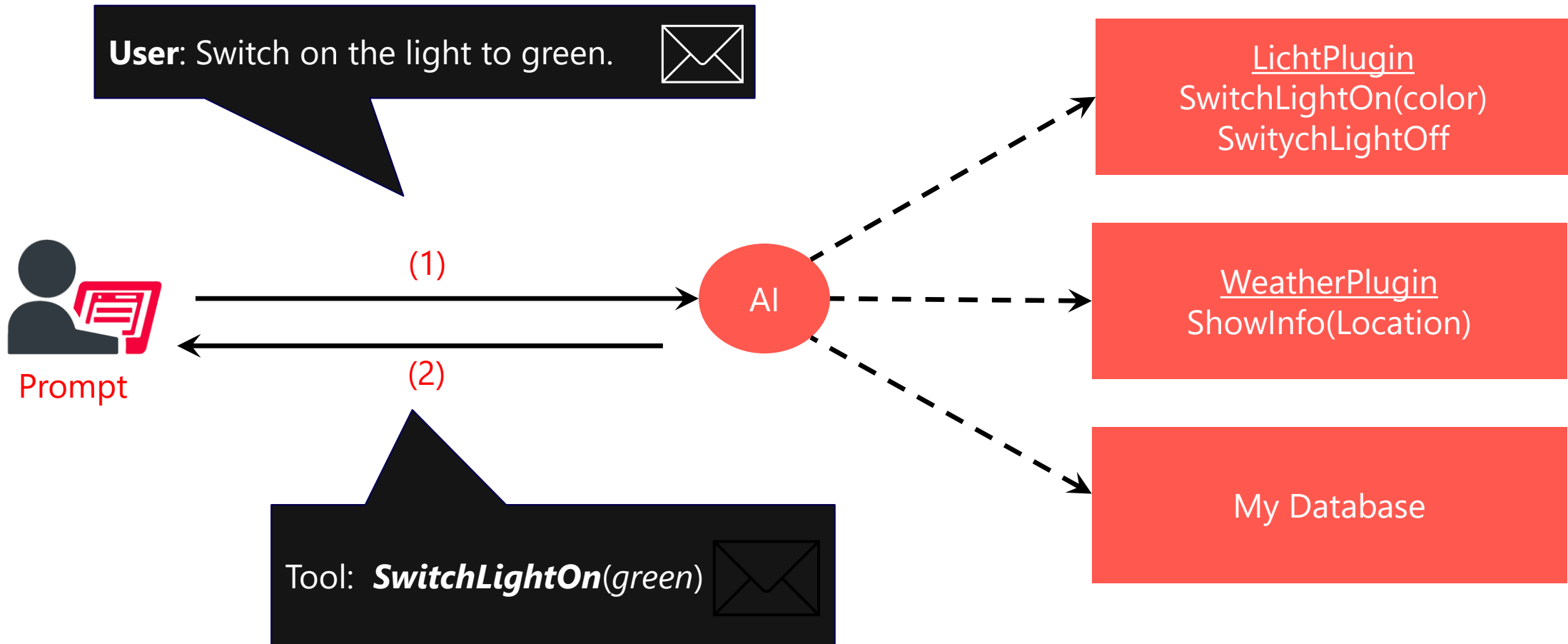
Function Calling



Chat Model Interaction



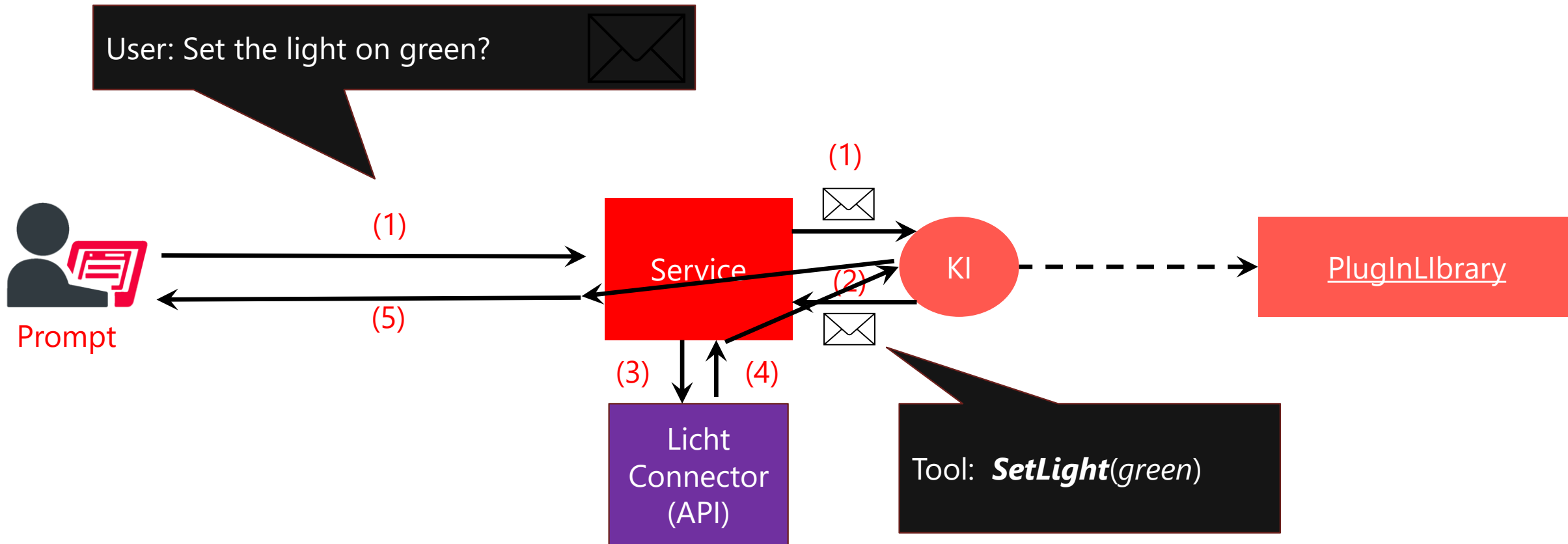
Function Calling



Function Calling

```
{
  "type": "function",
  "function": {
    "name": "get_weather",
    "description": "Retrieves current weather for the given location.",
    "strict": true,
    "parameters": {
      "type": "object",
      "properties": {
        "location": {
          "type": "string",
          "description": "City and country e.g. Bogotá, Colombia"
        },
        "units": {
          "type": ["string", "null"],
          "enum": ["celsius", "fahrenheit"],
          "description": "Units the temperature will be returned in."
        }
      },
      "required": ["location", "units"],
      "additionalProperties": false
    }
  }
}
```

Function Calling



GenAI unpacked: Beyond Basics



THANK YOU 😊

<https://www.linkedin.com/in/damirdobric/>

Dr. Damir Dobric

Lead Software Architect daenet GmbH / ACP Digital

Microsoft Regional Director,

Most Valuable Professional: AI

