



MARINA PANTCHEVA

DIRECTOR LINGUISTIC AI SERVICES

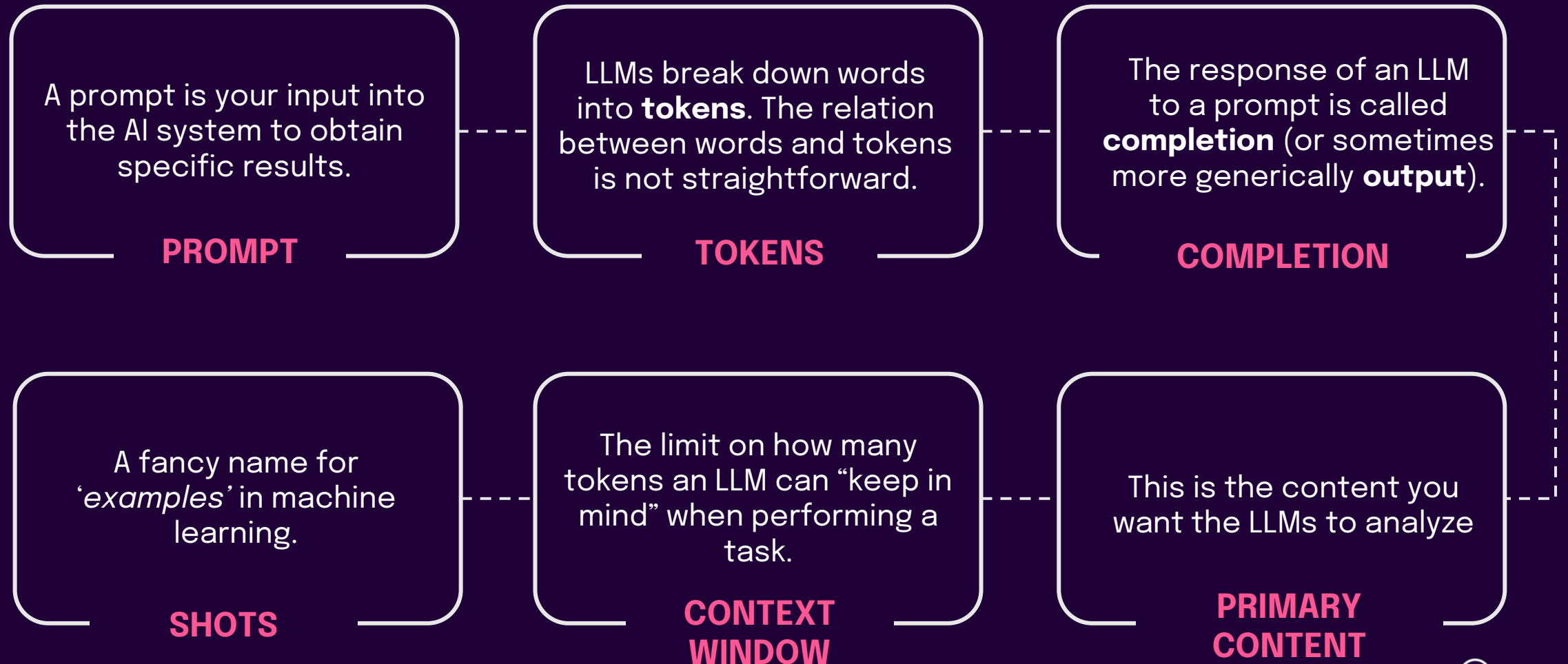
PROMPTING WITH PURPOSE

HANDS-ON ADVICE FOR A WIDE
AUDIENCE

GENAI IN LOCALIZATION CONFERENCE
MAY 28, 2025



6 BASIC CONCEPTS IN AI



LLM EVOLUTION

The great imitators

GPT 3, Llama 1, Mistral 7B



INFANCY

Instruction-tuned

*GPT 3.5, Llama 2,
Claude 3.5, Gemini 1.0*



CHILDHOOD

Agentic AI

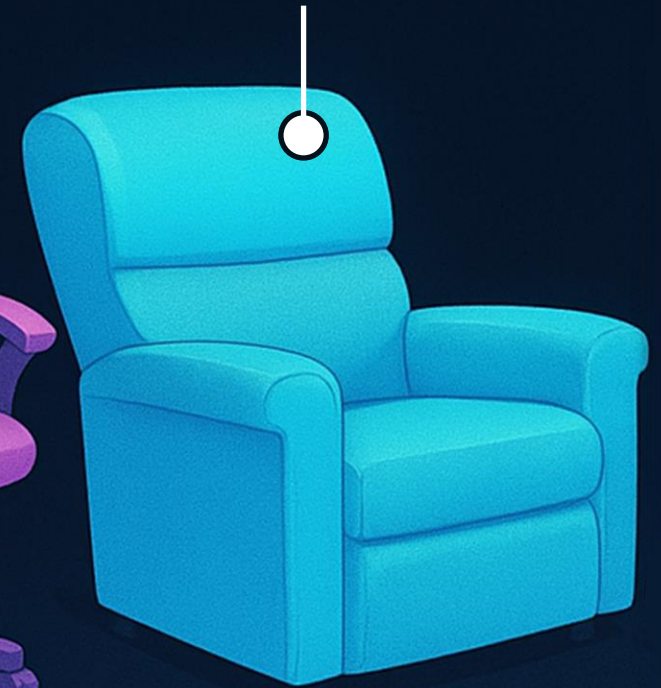
*AutoGPT, Operator,
Deep Research*



ADOLESCENCE

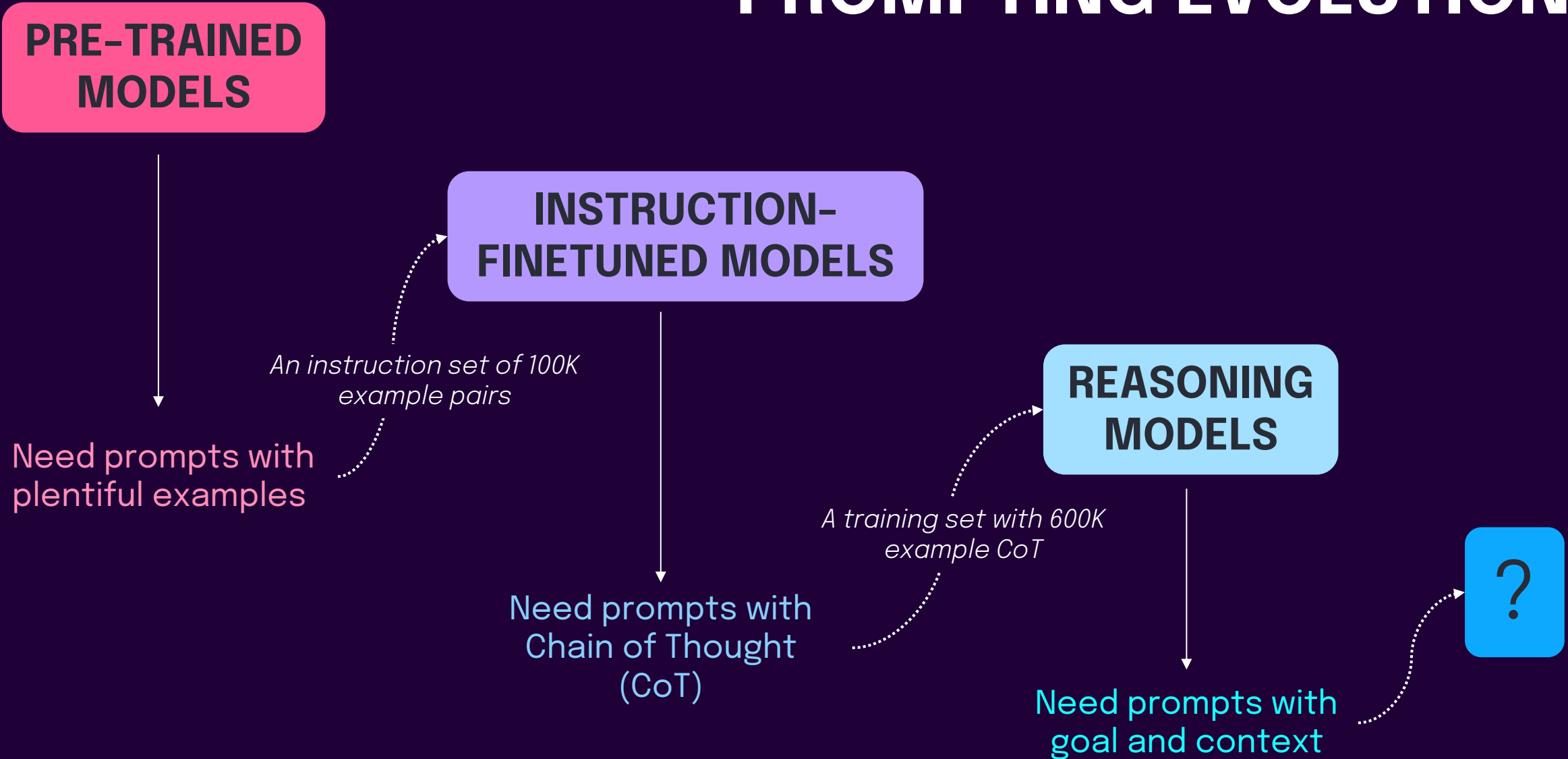
Reasoning models

*GPT o1/o3, Gemini Flash
DeepSeek R1/2*



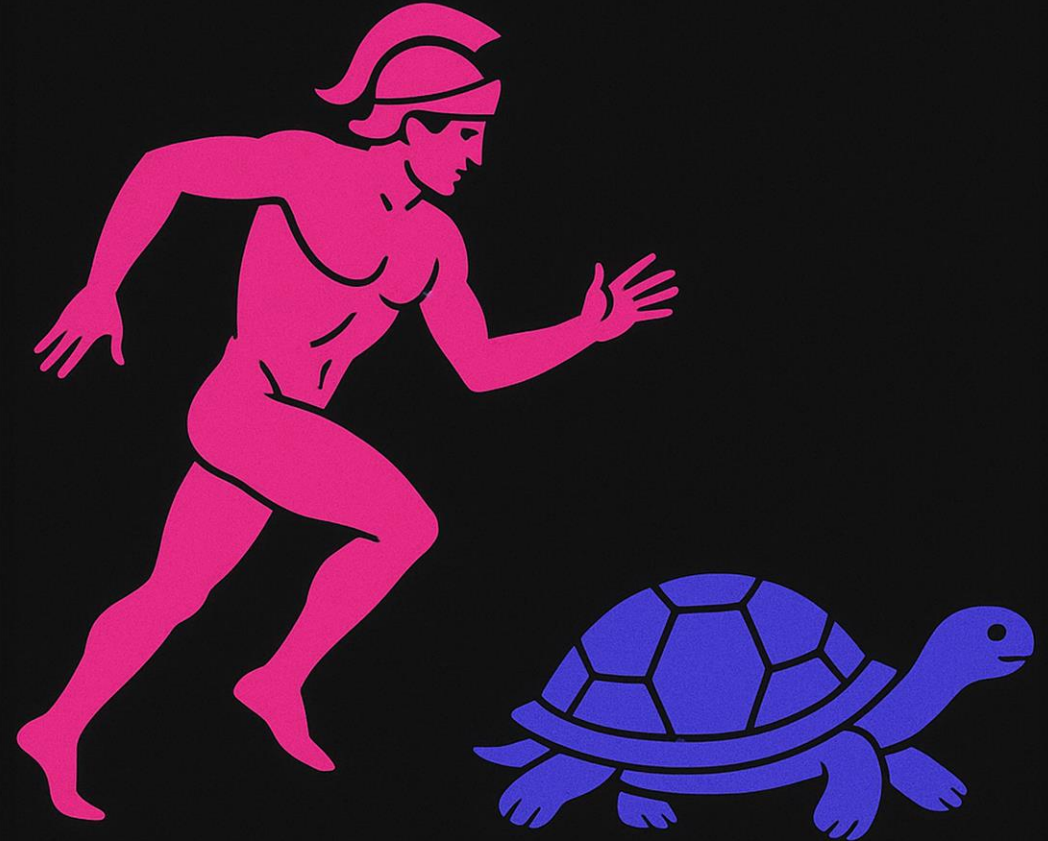
ADULTHOOD

PROMPTING EVOLUTION



ZENO's PARADOX OF CATCHING UP WITH AI

Chasing AI
developments only
halves the distance
to the **latest AI**
frontier model.



AI PILOTS

A detailed illustration of a rustic kitchen with warm lighting. On the left, a tall wooden shelf is filled with various jars and containers. In the center, a wooden table holds an open book, a vase of pink flowers, and several glass bottles. To the right, a wooden counter with a built-in oven and stovetop is visible, with a basket of red apples and other kitchen items. The background features a window with a warm glow and a hanging light fixture. Overlaid on the scene is a diagram with five white speech bubbles containing text, connected by lines to represent a workflow.

Cook
(prompt
designer)

Manual
tasting
(evaluation)

Tested
recipes
(prompts)

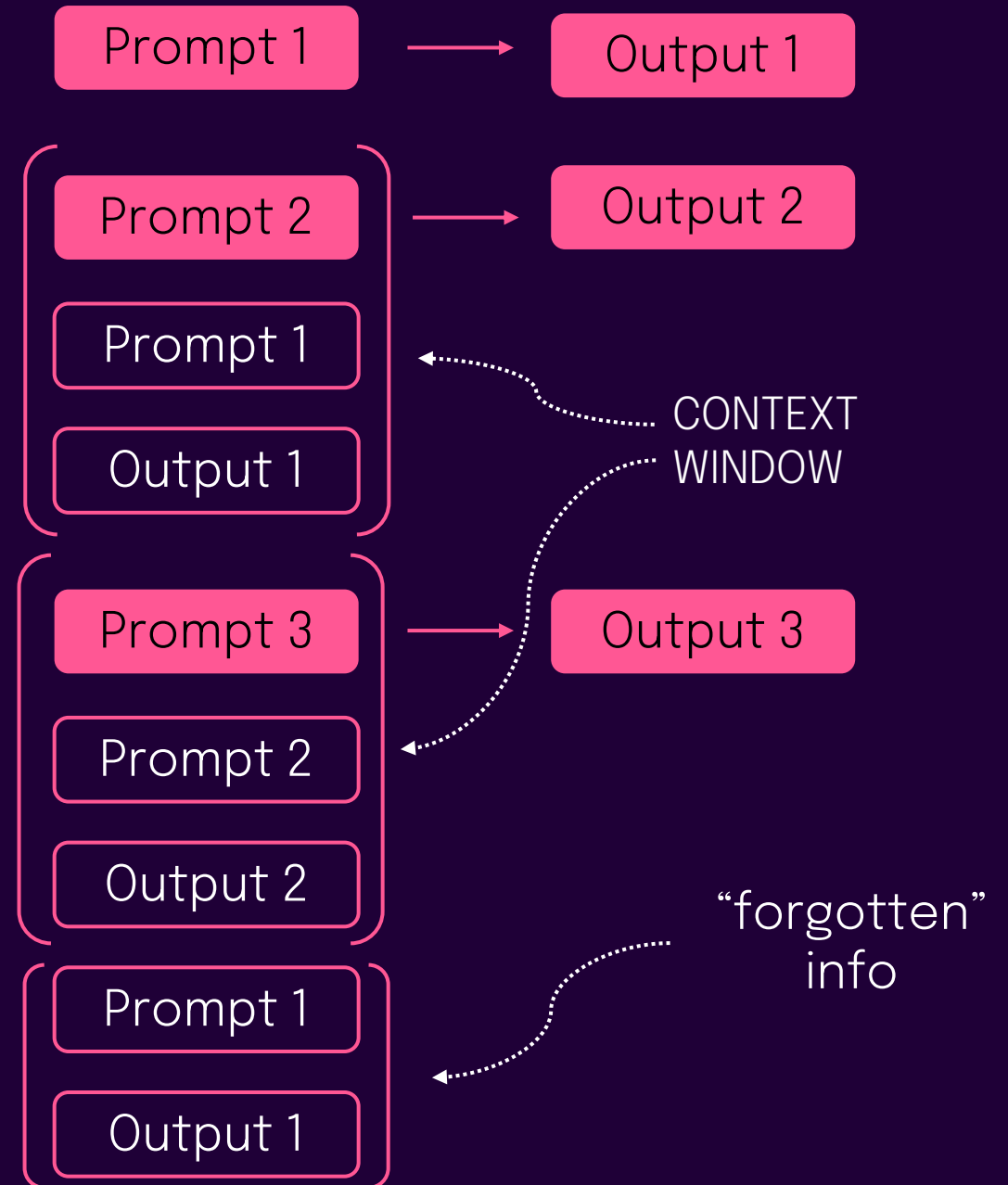
Selected
ingredients
(data)

Small
portions
(datasets)

CHATTING WITH AI

- Natural language dialogue
- **System prompt imposed:** No control over the actual prompt and parameters.
- **Contextual continuity:** the generated prompt contains all messages in the chat session that fit into the context window.
This incremental prompting may lead to degradation of response adequacy.

Brainstorming, question answering, assistance, quick proof of concepts



AI IN PRODUCTION

Many cooks
involved

Complex
infra-
structure

Giant
quantities
(TMs, TBs,
SGs)

Rare
ingredients
(low-resource
languages)

Multiple
orders
(continuous
production)



PROMPTING VIA API

- Technical knowledge needed to manage rate limits, authentication, etc.
- Full **control over the prompt** and its parameters
- **No contextual continuity**: context information is included in each call

Automated workflows, repeatable tasks, large scale operations, AI integration in apps

Prompt 1 → Output 1

Prompt 2 → Output 2

Prompt 3 → Output 3

Each API call (prompt session) runs in isolation.

THE TOKEN OVERHEAD PROBLEM

PROBLEM

The token count of the prompt outweighs by far the tokens of the text you want to process.

SOLUTION 1

Group the text segments into **batches** of 30–50.

- Wrap each segment in tags to avoid leakage.
- If concatenating, use unique delimiter that the LLM preserves.
- If using IDs, beware of ID shift.

SOLUTION 2

Extract into the prompt only the reference information that is **relevant** to the processed segments:

- relevant style guide rules
- terms that are indeed present in the segments

Query to an assistant

Prompting non-agents

Step-by-step instruction

Context and references

LLM output:

A **response** –
an answer or
converted input

Primary failure:

- wrong/irrelevant answers
- Incorrectly formatted/converted segments.

Desired output format

A few examples

Standard operating procedures

Prompting agents

Specification of goal and success

Catalogue of available tools

LLM output:

A **plan**
a JSON scratchpad

An **action**:
a tool call

A **decision**:
a termination signal

Primary failure:

- Wrong actions
- Overspending tokens

Method and memory

Guardrails and exit criteria

USE AI TO PROMPT AI

Have the LLM you plan to use for execution generate the prompt text.

This increases the chance that the prompt aligns with the instructions the LLM has seen in its training set.

Use the most advanced model available to draft the prompt. Reasoning models excel at this.

Ask the model you use for task execution produce a detailed prompt. Where it fails at instructing, it will also fail at executing.

PROMPT TROUBLESHOOTING

Break down the instructions into discrete steps within the same prompt.

1

Use prompt daisy-chaining where the output from prompt 1 is the input of prompt 2, etc.

2

Pull in more primary content to provide context (RAG).

3

Try a different LLM provider.

(The prompt may have to be adjusted)

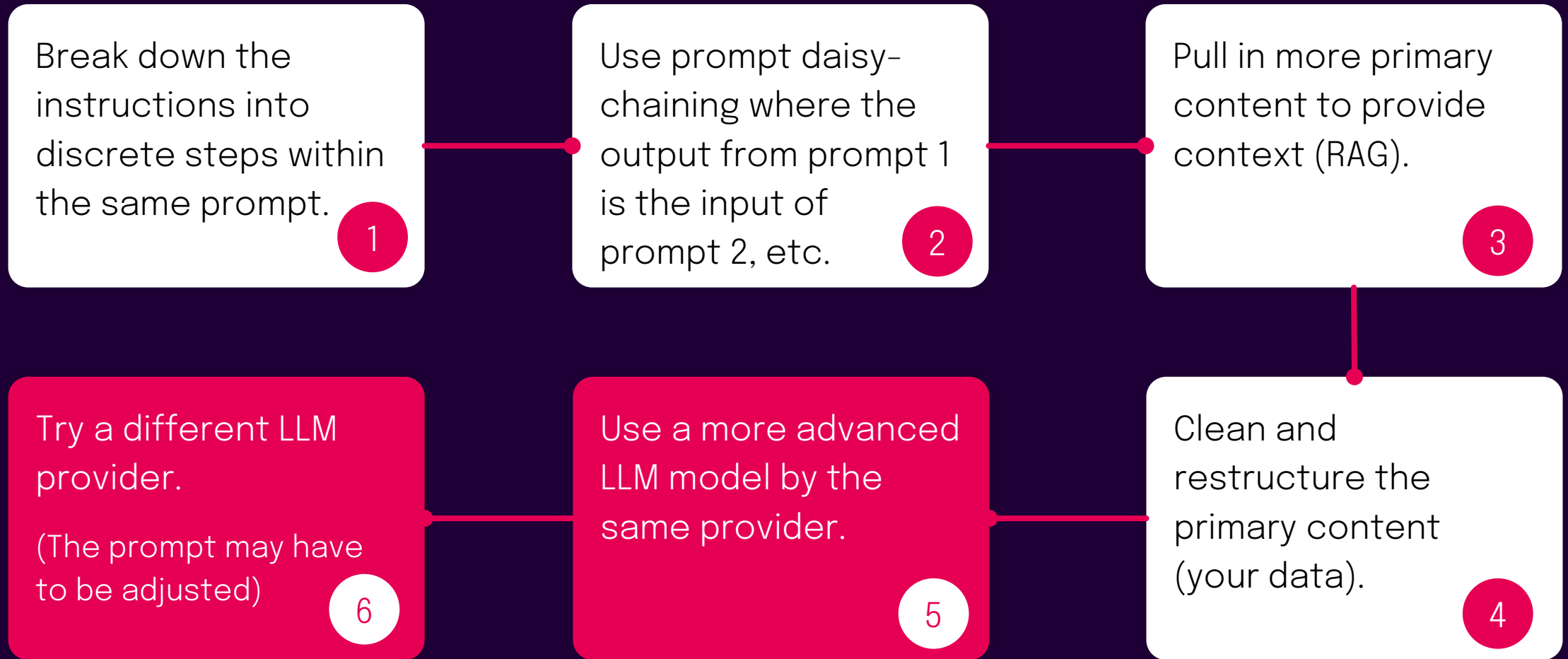
6

Use a more advanced LLM model by the same provider.

5

Clean and restructure the primary content (your data).

4



**Let's catch the AI
tortoise!**

