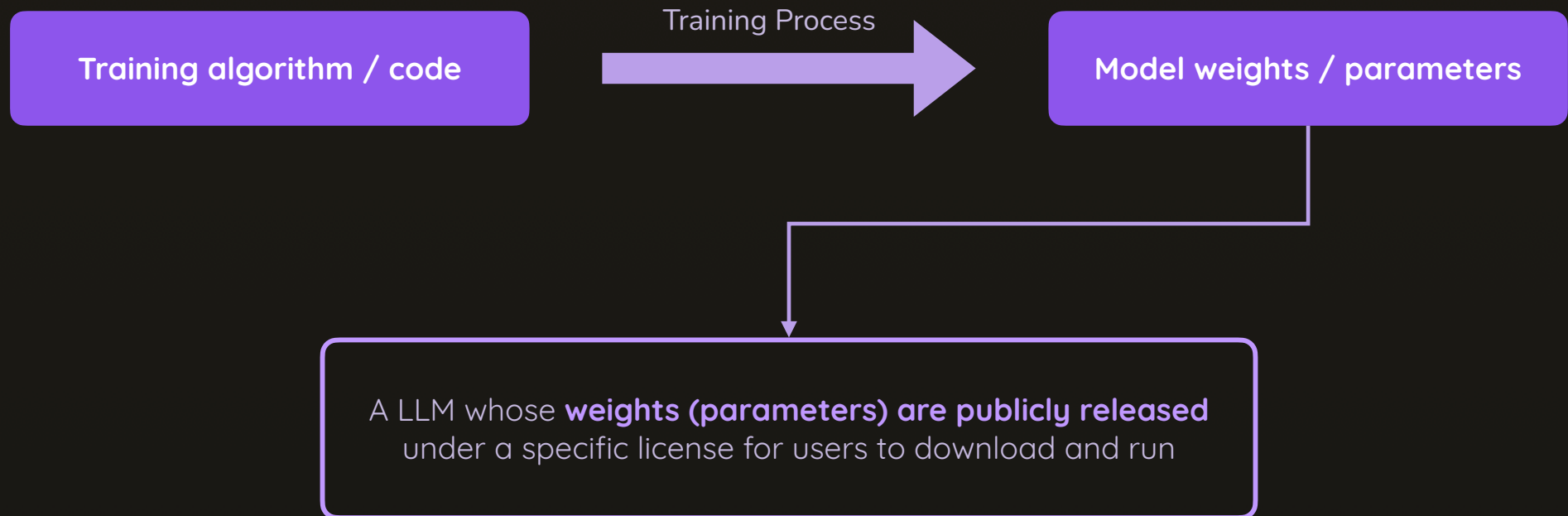


What Are “Open LLMs”?

What Are “Open LLMs”?



Open vs Proprietary LLMs

Open LLM

Free to use (must respect license, though)

Can be run locally or on your servers

100% privacy

No vendor lock-in, full control

Offline-first, low / no latency

Great for many (local) use-cases like
(internal) tools, text summarization, or
few-shot prompting

Proprietary LLM

Paid (Usage-based or Subscription)

Hosted by provider

Privacy depends on provider

Possible vendor lock-in, only little control

Internet connection required

Necessary for use-cases where you need
best-in-class performance

Examples for Open LLMs



Meta's Llama Models



Google's Gemma Models



DeepSeek Models

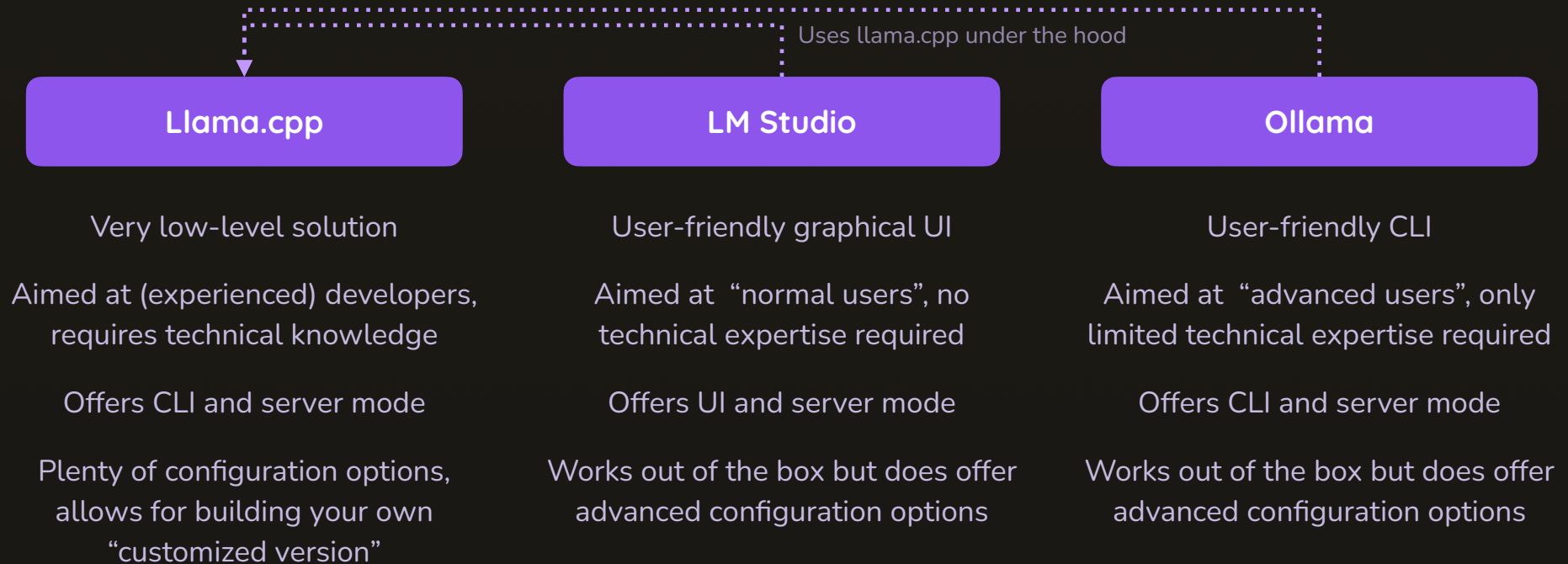


Some Mistral Models

Where Do You Find Open Models?

Open vs Proprietary LLMs

Options For Running Open LLMs Locally



Also: **On-demand usage** via **services like Groq**

Important: Licenses!

MIT / Apache 2.0

Very permissive licenses

Allow private & commercial use without limitations

No or only little attribution required

Llama

Somewhat permissive licenses

Allows private & commercial use **with** limitations

Attribution required

Prohibits usage for certain use-cases

Gemma

Somewhat permissive licenses

Allows private & commercial use **with** special requirements

Attribution required

Prohibits usage for certain use-cases

In general: Check the license & terms of the model you plan to use!

Beyond LLMs & Text Generation

We'll Run Models Locally
Not remotely



Hardware Requirements & Quantization

Running Large Models On “Small” Machines

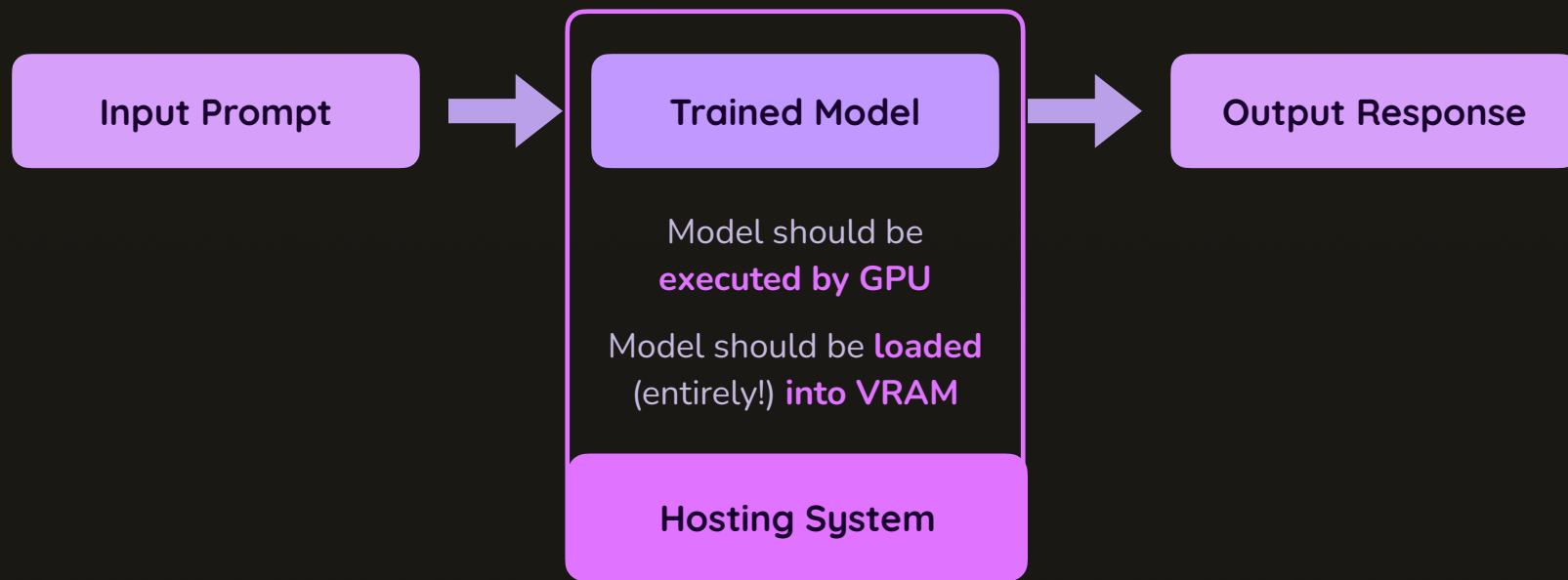
Understanding Model Parameters & Sizes

Hardware Requirements

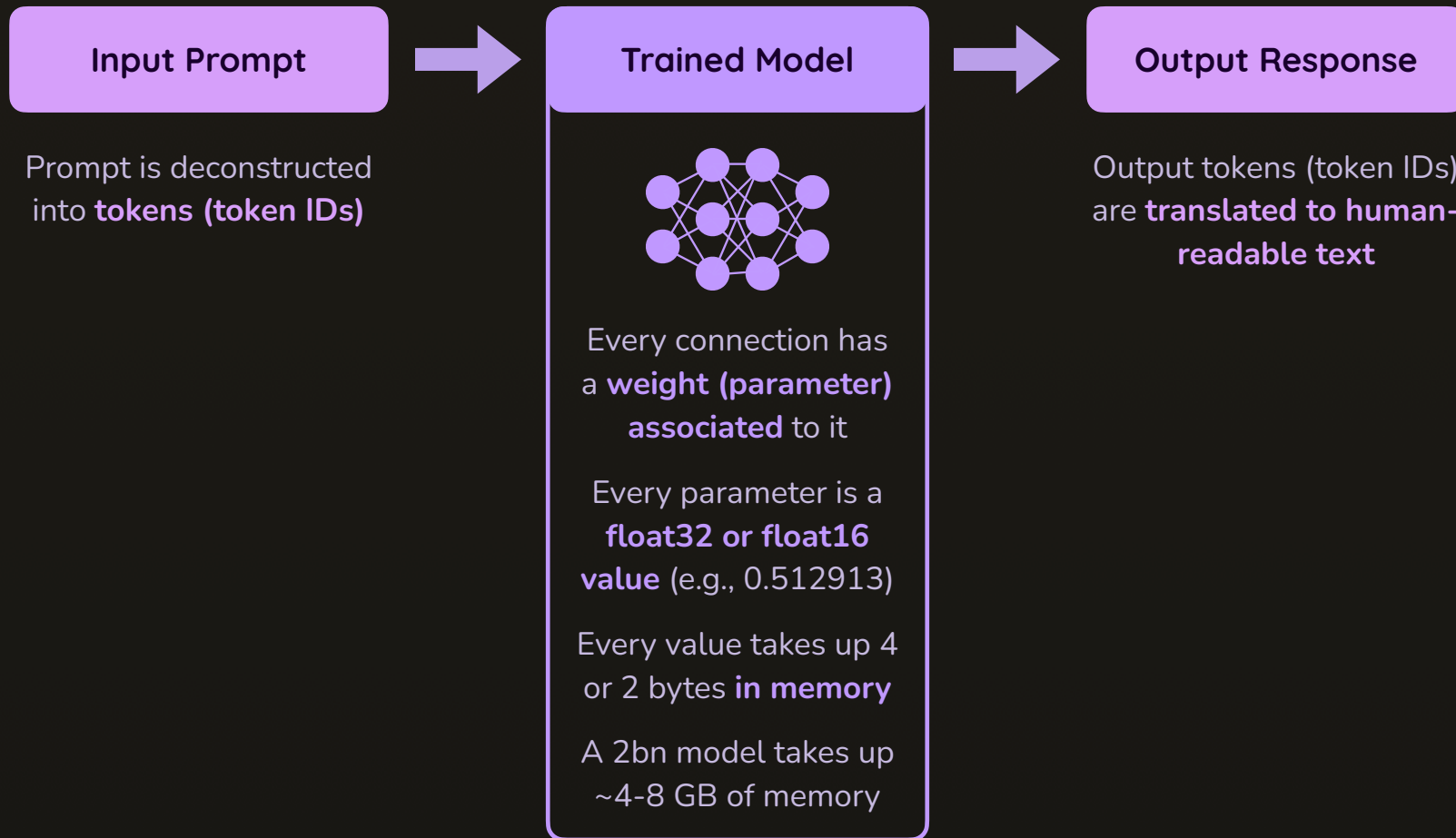
Exploring Quantization

Running LLMs For Inference

Inference is the process of using a (open or proprietary) LLM to generate output



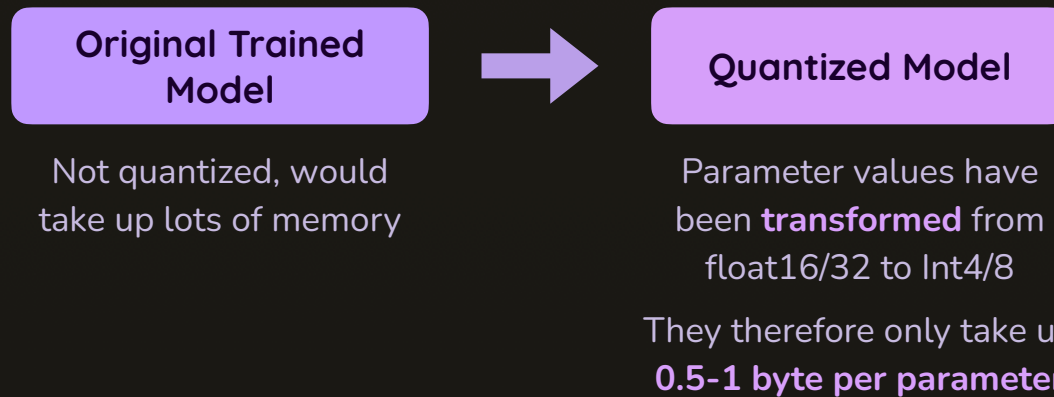
Understanding Model Parameters & Sizes



Does It Run On My Machine?

Understanding Quantization

In order to make LLMs **consume less (V)RAM** (and potentially also increase inference speed), the original models are **typically quantized (= compressed)**



You don't need to perform this quantization yourself!

Models shared on Huggingface, especially when usable via LM Studio or Ollama, are **available as quantized versions**



LM Studio

A Convenient Chat Interface For Running Open Models Locally

Installing & Using LM Studio

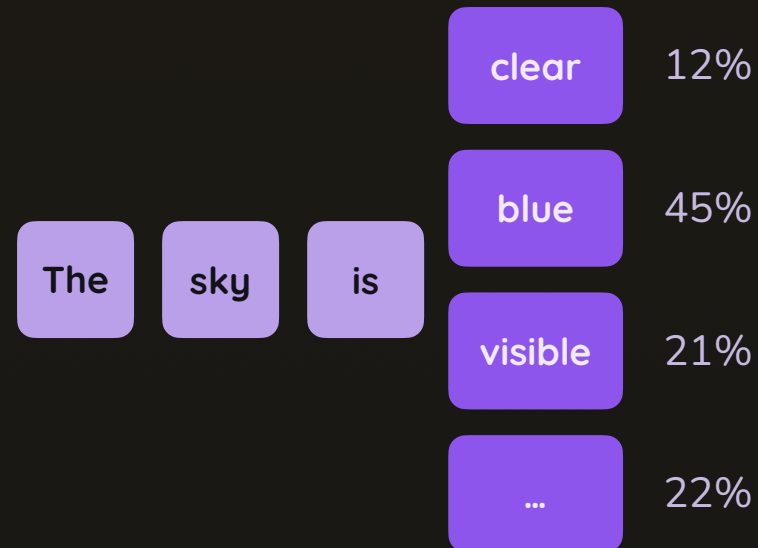
Advanced Configuration Options

Example Usage & Use-cases

Programmatic Usage

Tweaking “temperature”, “top_k”, “top_p” and “min_p”

LLMs generate token candidates



You can configure, which tokens are considered for sampling

Tweaking “temperature”, “top_k”, “top_p” and “min_p”

Temperature

Modifies probabilities

Low temperature:

Exaggerates differences

High temperature:

Flattens differences

top_k

Limits number of candidates

k = 5:

Choose from 5 most likely

k = 1:

Always choose the most likely

top_p

Limits candidates based on
combined probabilities

p = 0.5:

All candidates that combined have > 50%

p = 0.9:

All candidates that combined have > 90%

min_p

Discard candidates that don't
meet min-probability threshold

p = 0.05:

Remove candidates with probability < 5%

p = 0.5:

Remove candidates with probability < 50%



Ollama

A Minimal Solution For Advanced Users

Installing & Using Ollama

Advanced Configuration Options

Example Usage & Use-cases

Programmatic Usage



Llama.cpp

No Convenience, Full Control

Understanding llama.cpp

Installing & Using llama.cpp

Example Usage & Use-cases

Programmatic Usage