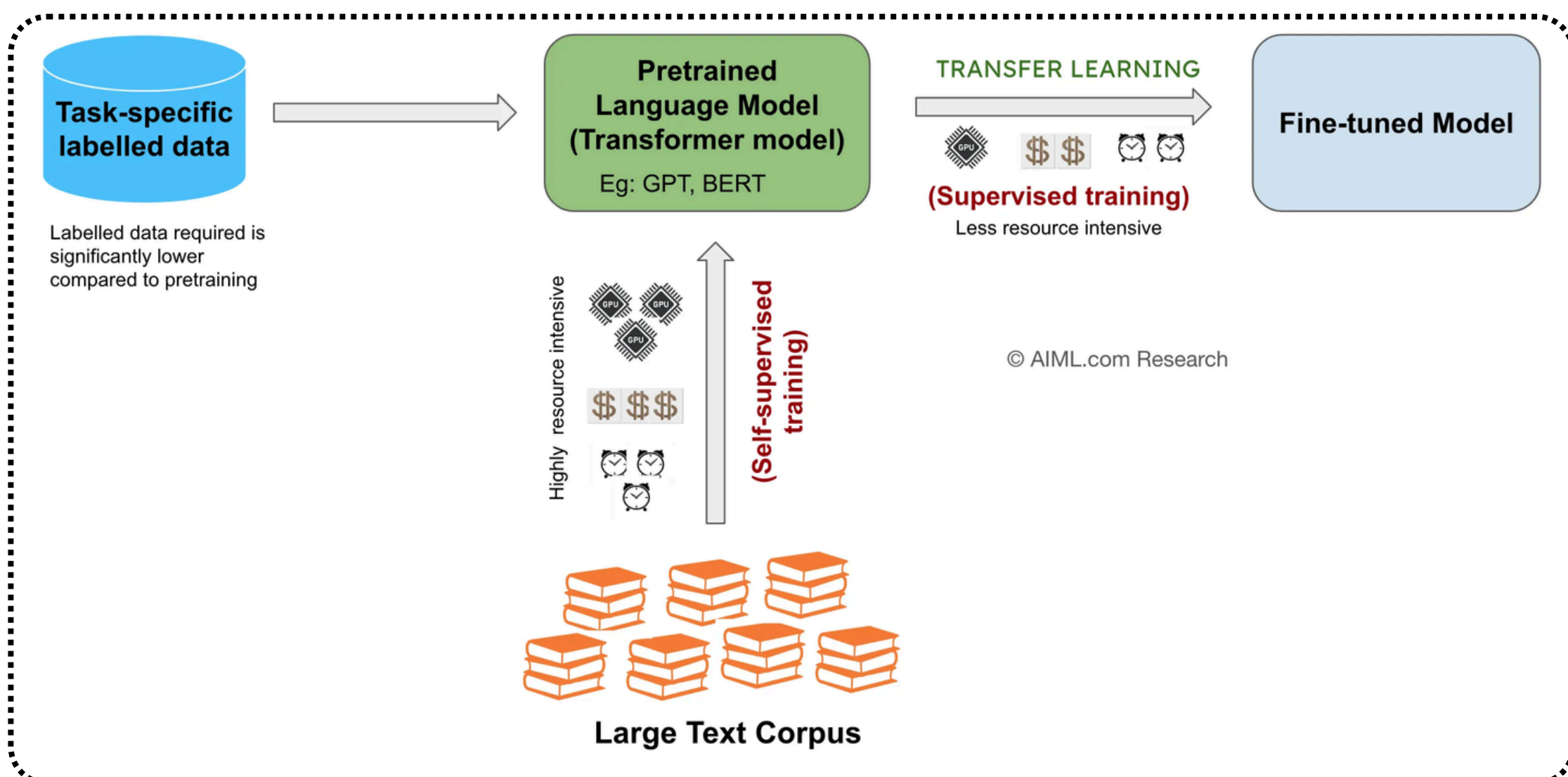
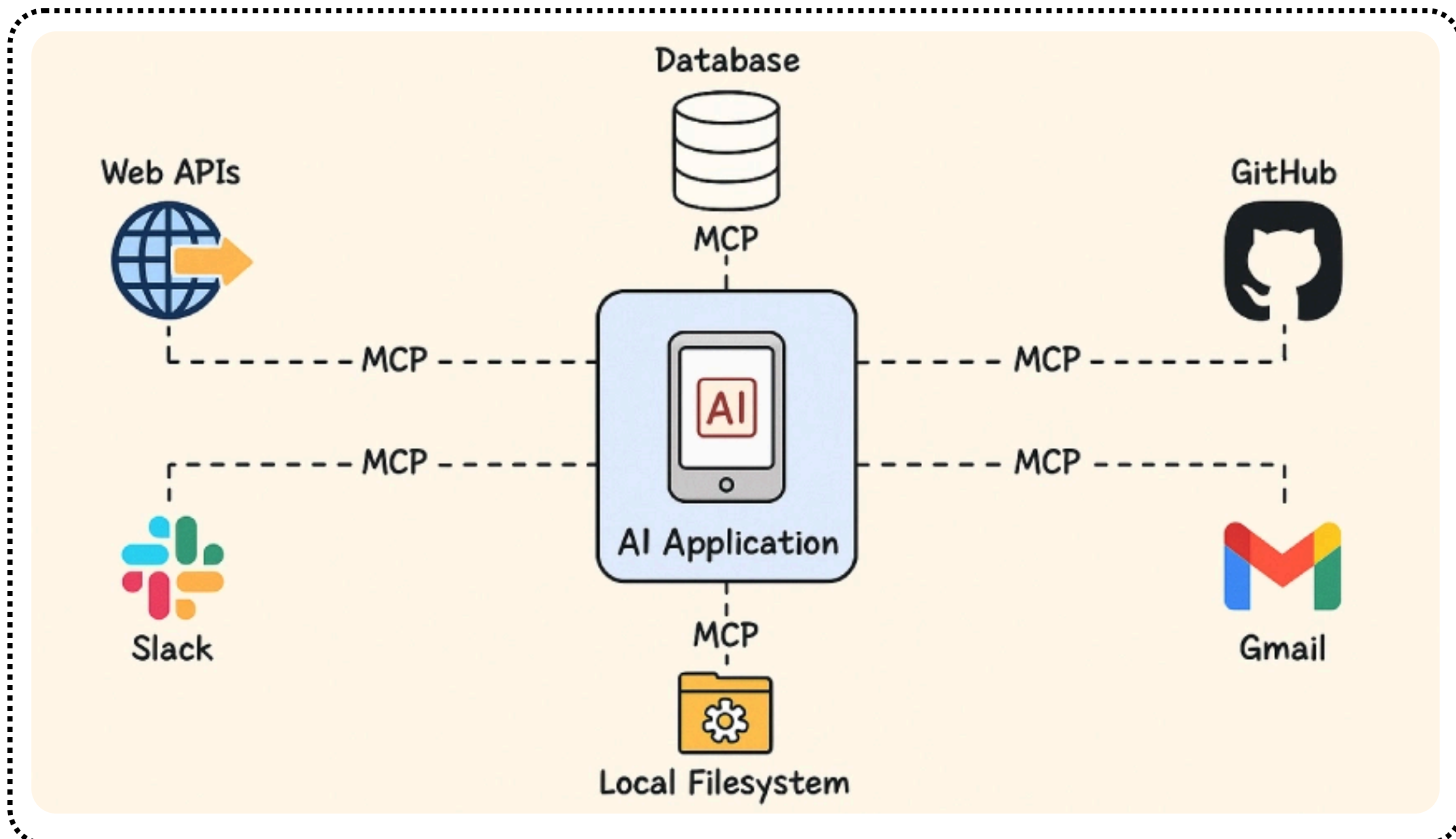


# MCP vs FineTuning



# What is Model Context Protocol?

MCP refers to injecting external information directly into the model's context window — like giving it background info, instructions, or examples at the time of inference.

## Examples include:

- Prompt engineering
- Few-shot or zero-shot learning
- RAG (Retrieval-Augmented Generation)
- Tool usage (e.g., APIs, code interpreters)

## Pros:

- No training needed – just update your prompts or context.
- Instant iteration – super fast to test changes.
- Dynamic & adaptive – great when info changes frequently (e.g., product catalogs, news).
- Lower cost – you don't need to re-train or host a custom model.

## Cons:

- Context window limit – models can only take in so much (e.g., 8K–128K tokens).
- Not sticky – the model “forgets” everything between queries.
- Prompt complexity – good prompting can be tricky and error-prone.

# What is FineTuning

Fine-tuning is about training the base model on your custom dataset, adjusting its internal weights to better reflect your domain or task.

## Pros:

- Permanent behavior change – no need to re-feed the same info.
- Efficient inference – shorter prompts once it's trained.
- Better task performance – especially in narrow or technical domains.
- Custom tone/style – e.g., matching a brand voice or legal language.

## Cons:

- Expensive to train and maintain – compute + storage + know-how.
- Slow iteration cycles – every tweak may need retraining.
- Outdated fast – static knowledge can quickly go stale.
- Risk of overfitting – especially with small or biased datasets.

## Best For:

- Specific tasks (e.g., classification, summarization)
- Domains with stable, long-term knowledge
- Custom agents with consistent tone or logic

# Difference between MCP & FineTuning

Feature	Model Context Protocol	Fine-Tuning
Setup Time	Minimal	High
Update Flexibility	Very high	Low
Cost	Low to moderate	High
Performance	Good (depends on prompt)	Excellent (task-specific)
Knowledge Permanence	Temporary (per session)	Permanent (until re-trained)
Best For	Dynamic, user-specific tasks	Stable, high-accuracy use cases

## Which Should You Use?

Here's a quick suggestion framework:

If you need...	Go with...
Quick iteration & low overhead	Model Context Protocol
To serve up-to-date info	Context / RAG
A model that "remembers" your domain	Fine-Tuning
A specific tone or behavior across all replies	Fine-Tuning
Personalized experience per user	Context / Tools / Functions
Long-term investment for scale	Fine-Tuning (with eval/monitoring)