

Structured outputs

LLMs with tools

- Allow LLM to output tool calls
- Special tags
- Special chat-template



Structured output

- What if we only want to **parse** output of LLM?
- Option 1: In context example

DEMO



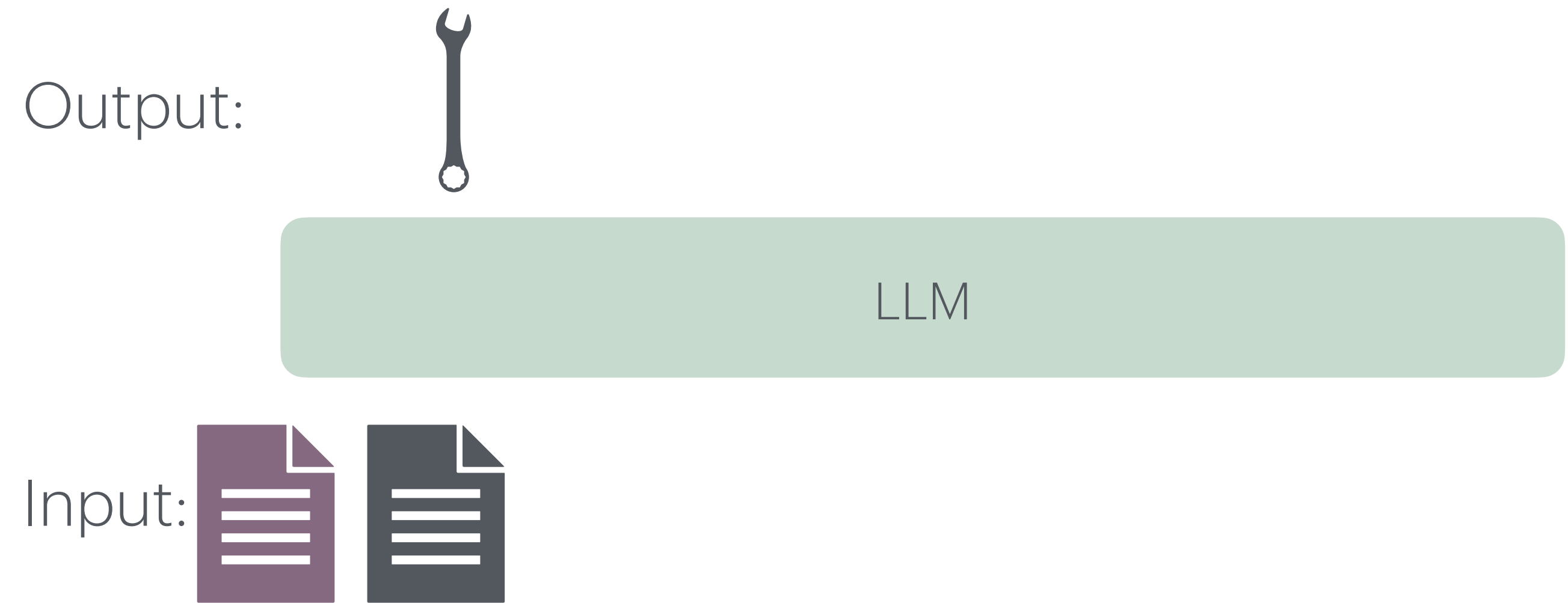
Structured output

- What if we only want to **parse** output of LLM?
- Option 1: In context example
 - Parsing can easily fail (more later)



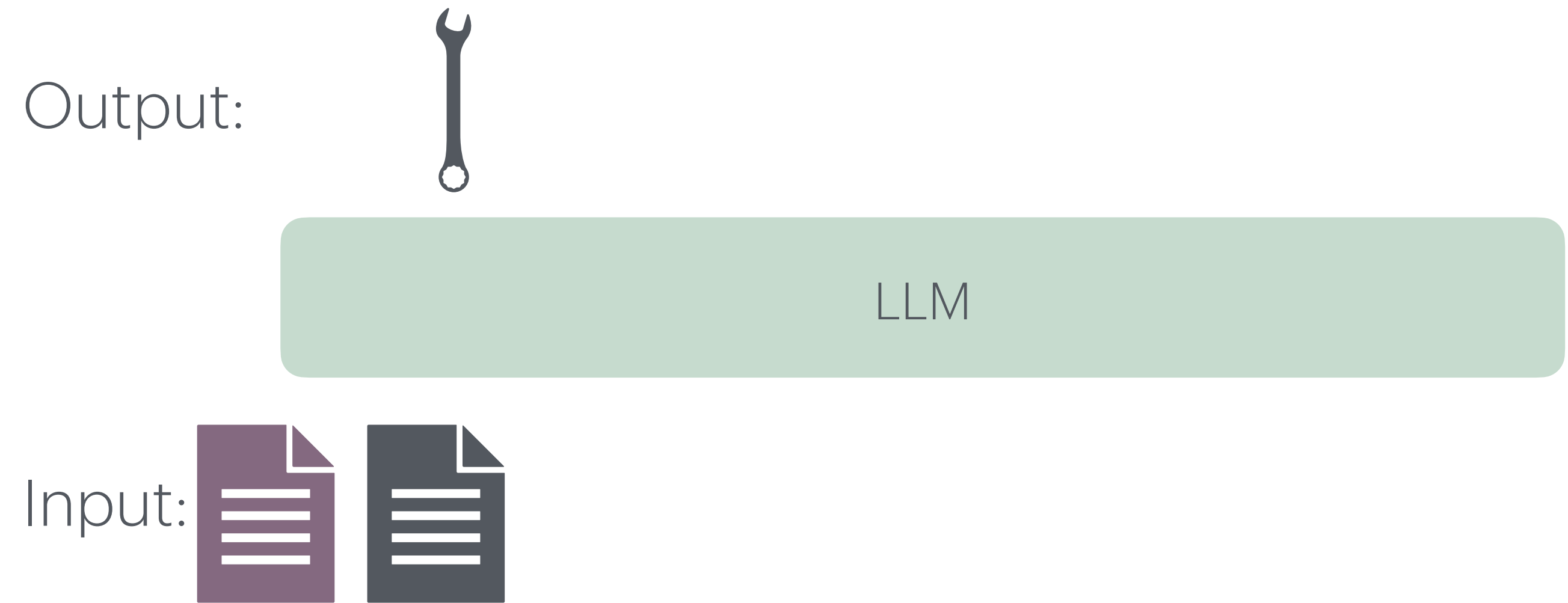
Structured output

- What if we only want to **parse** output of LLM?
- Option 1: In context example
 - Parsing can easily fail (more later)
- Option 2: Use a tool, arguments = json fields



Structured output

- What if we only want to **parse** output of LLM?
 - Option 1: In context example
 - Parsing can easily fail (more later)
 - Option 2: Use a tool, arguments = json fields
 - More training data
 - Parsing might still fail



Structured output

- What if we only want to **parse** output of LLM?
 - Option 1.1
 - Write a robust parser (in python)
 - Let LLM know that you failed to parse
 - Hope for the best
 - Option 1.2: Constrain Decoding

