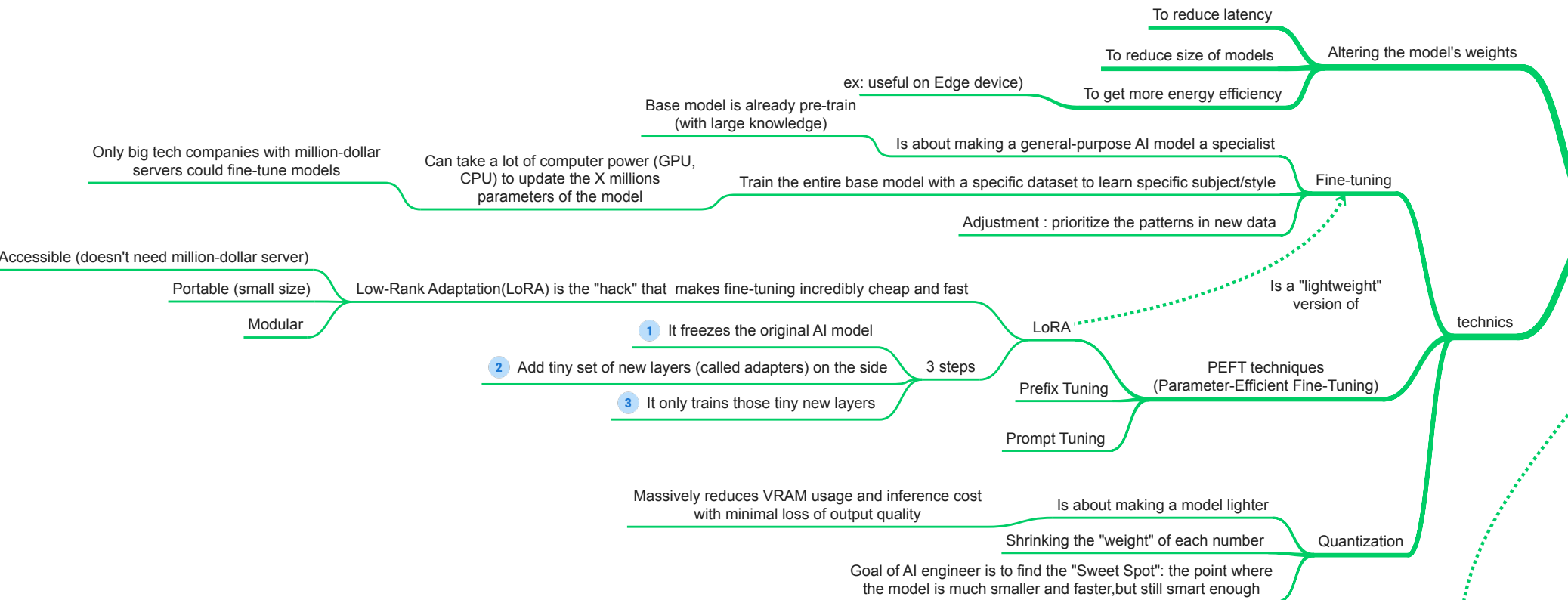
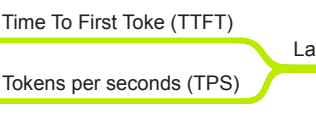


**Google Cloud GenAI - Model Optimization**  
<https://squasta.github.io>  
Updated 13 April 2026

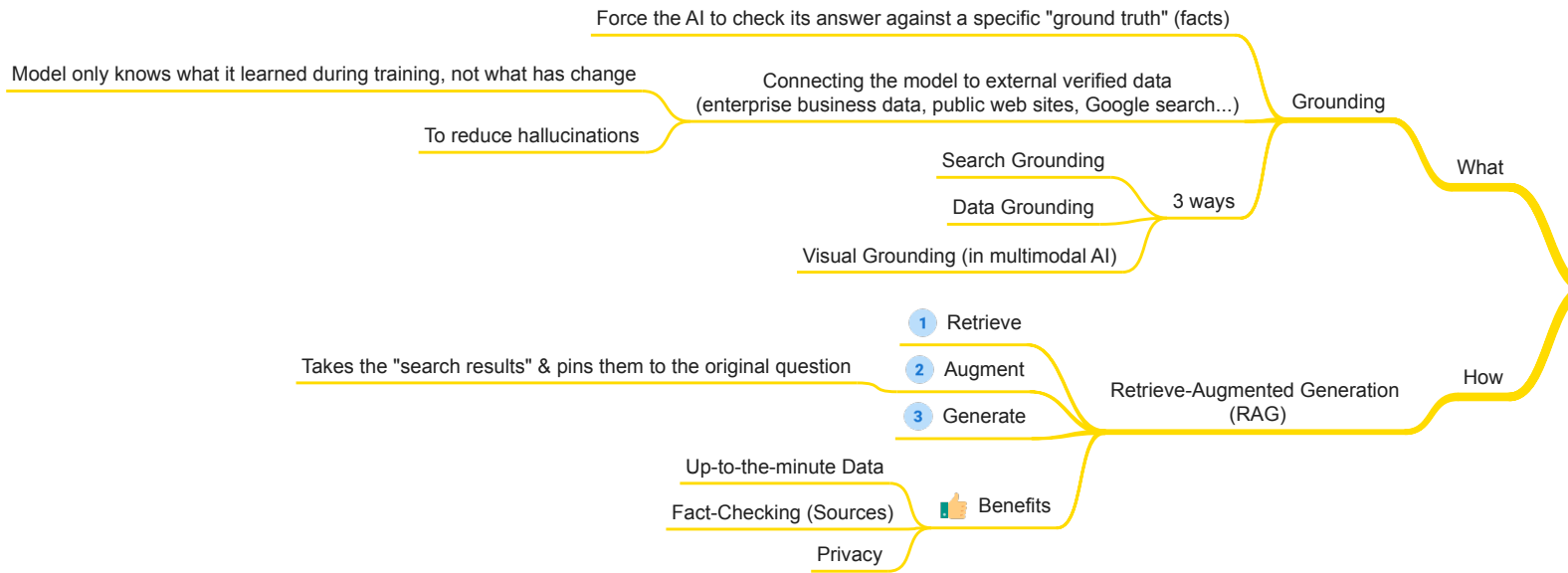
**Model Adaptation (High cost / effort)**



**Optimizing infrastructure**

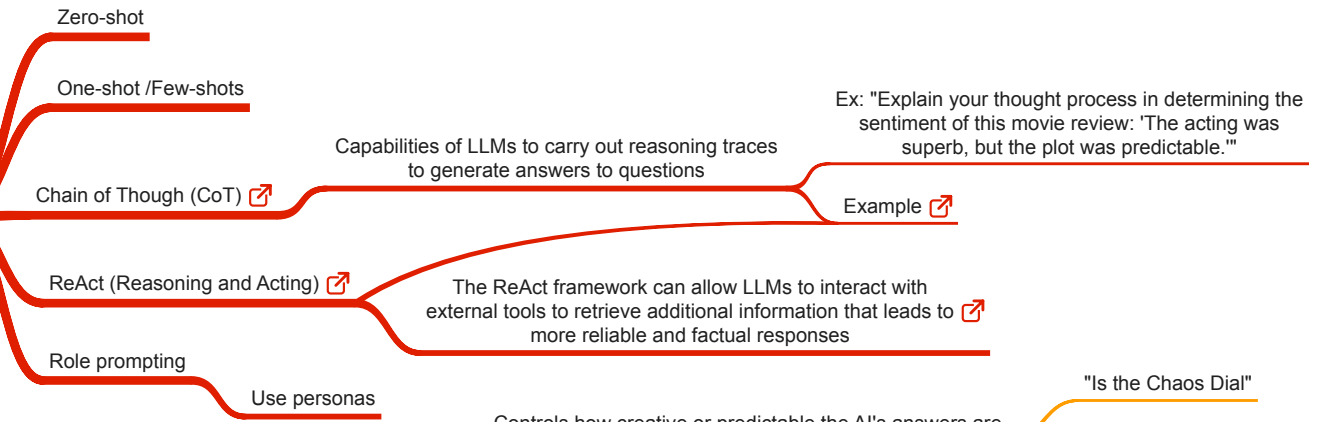


**Grounding & RAG (Medium cost / effort)**

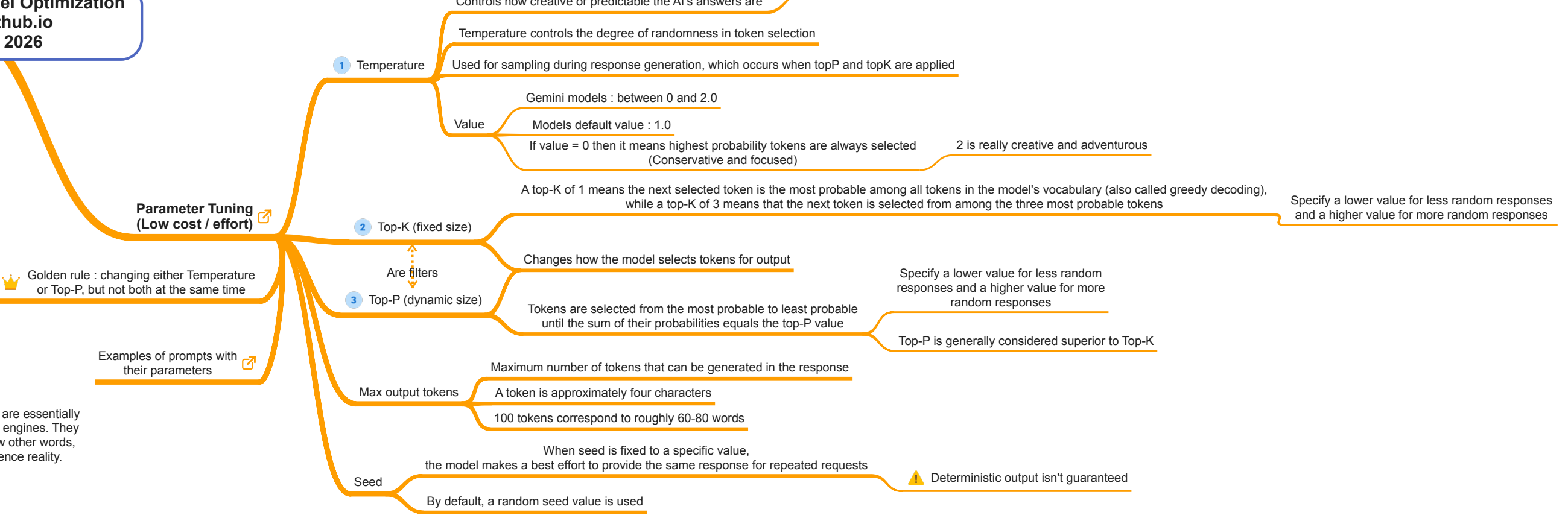


Large Language Models (LLMs) are essentially super-advanced "autocomplete" engines. They know which words usually follow other words, but they don't actually experience reality.

**Prompt Engineering (Low cost / effort)**



**Parameter Tuning (Low cost / effort)**



Examples of prompts with their parameters