CS 696 Applied Large Language Models
Spring Semester, 2025
Doc 26 End Comments
May 1, 2025

# Projects

What did you use - in detail

Code

Source code

All dependencies with enough information for me to install

Pip commands, requirements.txt, and which versions

Data

Which data sources?

How do I access the data

References

Websites, articles, books, and AI you used

How did you use them

Where in your project

Hardware used

# Projects

What did you do

Each of you is doing something different

Don't make me read source code to figure out what you are doing

What was the goal of your project?

What did you do to achieve your goal

# Projects

What were the results?

Don't make me interpret your output to figure out the outcome

What does your output tell you, and why

Did you achieve your goal

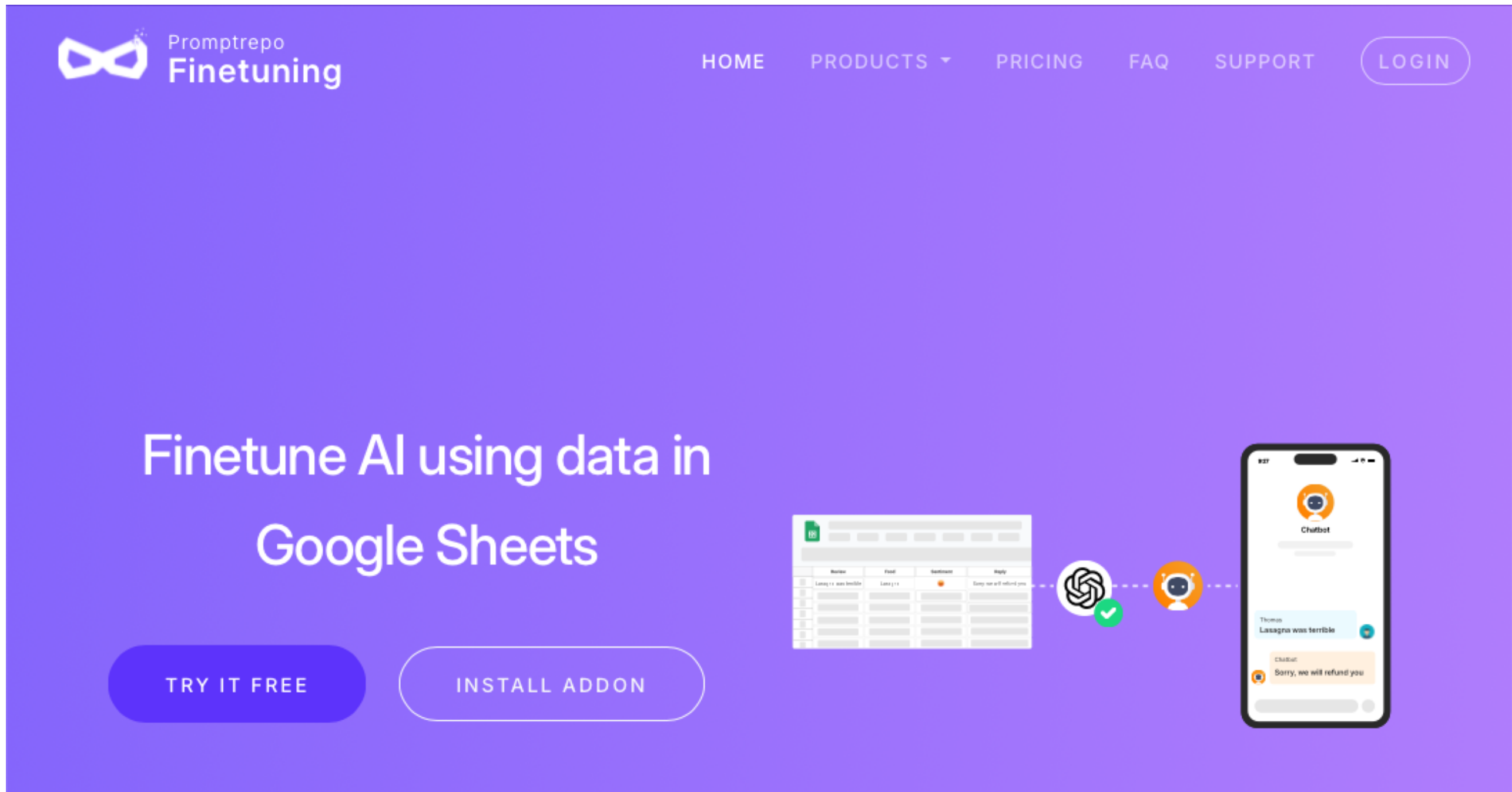# Projects

Known issues and limitations

Known bugs

Things that don't work

# Projects

Don't make me search for all of the above

# News



https://promptrepo.com/finetune/

# Finetuning

# Defeating Prompt Injections by Design

# Defeating Prompt Injections

# Defeating Prompt Injections by Design

# Model Released This Week (So far)

DeepSeek-Prover-V2

# Model Released This Week (So far)

Amazon Nova Premier

Our most capable model for complex tasks and teacher for model distillation

| | | Nova Pro | Nova Premier |
|---|---|---|---|
| **Text intelligence** | Undergraduate level knowledge<br>MMLU | 85.9% | 87.4% |
| | Science<br>GPQA Diamond | 50.0% | 57.1% |
| | High school math competition<br>AIME 2025 | 5.3% | 16.0% |
| | Math problem-solving<br>MATH-500 | 76.6% | 82.0% |
| | Coding<br>BigCodeBench Hard | 22.3% | 28.1% |
| | Coding<br>MBXP (5 languages) | 65.9% | 78.4% |
| | Instruction Following<br>IFEval | 92.1% | 91.5% |
| **Visual intelligence** | Visual understanding<br>MMMU | 62.0% | 68.0% |
| | Document understanding<br>OCRBench-v2 | 53.7% | 56.9% |
| | Chart understanding<br>CharXiv (Descriptive/Reasoning) | 70.5%/<br>40.6% | 84.6%/<br>48.8% |
| | Long-form video language understanding<br>EgoSchema | 72.1% | 73.8% |
| | Visual counting<br>TallyQA | 54.0% | 61.5% |
| **flows** | Retrieval-augmented generation<br>SimpleQA (SerpAPI) | 84.6% | 86.3% |
| | Function calling<br>BFCL (2025-04-25) | 60.8% | 63.7% |

# Model Released This Week (So far)

Phi 4

Reasoning,

14 B parameter,

Fine-tuned with reasoning demonstrations from OpenAI 03-mini

Reasoning-plus

Further trained with RL

Mini-Reasoning

Fine-tuned with synthetic data generated by Deepseek-R1



Phi-4 offers high quality results at a small size

# Model Released This Week (So far)

Qwen 3

32B, 14B, 4B, 1.7B, 0,6B

Hybrid thinking

| | Qwen3-235B-A22B | Qwen3-32B | OpenAI-o1 | Deepseek-R1 | Grok 3 Beta | Gemini2.5-Pro | OpenAI-o3-mini |
|---|---|---|---|---|---|---|---|
| | MoE | Dense | 2024-12-17 | | Think | | Medium |
| ArenaHard | 95.6 | 93.8 | 92.1 | 93.2 | - | 96.4 | 89.0 |
| AIME'24 | 85.7 | 81.4 | 74.3 | 79.8 | 83.9 | 92.0 | 79.6 |
| AIME'25 | 81.5 | 72.9 | 79.2 | 70.0 | 77.3 | 86.7 | 74.8 |
| LiveCodeBench v5, 2024.10-2025.02 | 70.7 | 65.7 | 63.9 | 64.3 | 70.6 | 70.4 | 66.3 |
| CodeForces Elo Rating | 2056 | 1977 | 1891 | 2029 | - | 2001 | 2036 |
| Aider Pass@2 | 61.8 | 50.2 | 61.7 | 56.9 | 53.3 | 72.9 | 53.8 |
| LiveBench 2024-11-25 | 77.1 | 74.9 | 75.7 | 71.6 | - | 82.4 | 70.0 |
| BFCL v3 | 70.8 | 70.3 | 67.8 | 56.9 | - | 62.9 | 64.6 |
| MultiIF 8 Languages | 71.9 | 73.0 | 48.8 | 67.7 | - | 77.8 | 48.4 |

1. AIME 24/25: We sample 64 times for each query and report the average of the accuracy. AIME'25 consists of Part I and Part II, with a total of 30 questions.
2. Aider: We didn't activate the think mode of Qwen3 to balance efficiency and effectiveness.
3. BFCL: The Qwen3 models are evaluated using the FC format, while the baseline models are assessed using the highest scores obtained from either the FC or prompt formats.

# Model Released This Week (So far)

Claude Integrations

      Claude works with desktop apps and remote servers

      Uses  Model Context Protocol (MCP)

Claude's Research

      Can search
         Web
         Google Workspace
         Integrations

    MCP documentation
      https://modelcontextprotocol.io/introduction

# Stable Diffusion

# Diffusion

Training is different than text LLM

Forward Process (Diffusion)
Incrementally add Gaussian noise until the image is pure noise

Reverse (Generative) Process
Image can be covered by removing the added noise step by step

Network is trained by
Take the noisy data $x_t$ at a particular step t as input

Predict the noise ($\epsilon$) that was added to get from $x_{t-1}$ to $x_t$

The model then generates images by
Starting with noise
Incrementally remove the noise

# Diffusion

High-Quality Generation

Stable Training

Slow Sampling

Conditioning
Trained to use text inputs

# Model Released This Week (So far)

Mercury from Inception

Trained by diffusion

# Model Released This Week (So far)

Mercury from Inception

# Mercury Coder

Try our first commercial-grade diffusion LLM

Write a simulator for 5 balls bouncing on a billiard table. Make collision physics realistic, without gravity. Use Javascript.scree

⚡ Suggested

**Write a simulator for 5 balls bouncing**
on a billiard table

**Make a particle system**
where particles follow the mouse cursor

**Illustrate a forward diffusion process**
in HTML 5

By using Mercury Coder, you agree to our Terms of Service and have read our Privacy Policy.

| | Mercury Coder Mini | Mercury Coder Small | Gemini 2.0 Flash-Lite | Claude 3.5 Haiku | GPT-4o Mini | Qwen 2.5 Coder 7B | DeepSeek Coder V2 Lite |
|---|---|---|---|---|---|---|---|
| **HumanEval** | 88.0 | 90.0 | 90.0 | 86.0 | 88.0 | 90.0 | 92.1 |
| **MBPP** | 77.1 | 76.6 | 75.0 | 78.0 | 74.6 | 80.0 | 81.0 |
| **EvalPlus** | 78.6 | 80.4 | 77.3 | 75.1 | 78.5 | 79.3 | 82.1 |
| **MultiPL-E** | 74.1 | 76.2 | 79.5 | 72.3 | 72.0 | 75.3 | 79.1 |
| **LiveCodeBench** | 17.0 | 25.0 | 18.0 | 31.0 | 23.0 | 9.0 | 37.8 |
| **BigCodeBench** | 42.0 | 45.5 | 44.4 | 45.4 | 46.8 | 41.4 | 50.0 |
| **Fill-in-the-Middle** | 82.2 | 84.8 | 60.1 | 45.5 | 60.9 | 56.1 | 46.9 |

# Programming & AI

Continue

Cline

Roo Code

Cursor

Vibe Programming

# Programming & AI

Two distinct patterns

Bootstrappers
  Tools: Bolt, v0

    Start with a design or rough concept
    Use AI to generate a complete initial codebase
    Get a working prototype in hours or days instead of weeks
    Focus on rapid validation and iteration

Iterators
  Tools: Cursor, Cline, Copilot, and WindSurf

    Using AI for code completion and suggestions
    Leveraging AI for complex refactoring tasks
    Generating tests and documentation
    Using AI as a "pair programmer" for problem-solving

# Programming & AI - Common Failure Patterns

# Programming & AI - Common Failure Patterns

The Demo-Quality Trap

AI make it easy to develop demo-quality software

But
Not complete

Hard to understand

Difficult to modify

# The Golden Rules of AI Coding

Be specific and clear about what you want

Always validate AI output against your intent

Treat AI as a junior developer (with supervision)

Use AI to expand your capabilities, not replace your thinking

Coordinate upfront among the team before generating code

Treat AI usage as a normal part of the development conversation

Isolate AI changes in Git by doing separate commits

Ensure that all code, whether human or AI-written, undergoes code review

Don't merge code you don't understand

Prioritize documentation, comments, and ADRs

Share and reuse effective prompts

Regularly reflect and iterate

# Staying Current

Simon Willison Blog

    https://localforge.dev/blog

  O'Reilly Online

    https://learning.oreilly.com/

https://medium.com

    Good for learning what exists

Vibe Coding: The Future of Programming

Addy Osmani

O'Reilly Media, Inc., **August** 2025

# Hacker News

https://news.ycombinator.com

Y **Hacker News**  new | past | comments | ask | show | jobs | submit

1. ▲ Linux Kernel Exploitation: Attack of the Vsock (hoefler.dev)
   89 points by todsacerdoti 3 hours ago | hide | 24 comments

2. ▲ Mercury, the first commercial-scale diffusion language model (inceptionlabs.ai)
   17 points by HyprMusic 30 minutes ago | hide | discuss

3. ▲ Zhaoxin's KX-7000 (chipsandcheese.com)
   35 points by ryandotsmith 1 hour ago | hide | 7 comments

4. ▲ Reversible computing with mechanical links and pivots (tennysontbardwell.com)
   90 points by tennysont 4 hours ago | hide | 37 comments

5. ▲ NotebookLM Audio Overviews are now available in over 50 languages (blog.google)
   195 points by saikatsg 4 hours ago | hide | 55 comments

6. ▲ Xiaomi MiMo Reasoning Model (github.com/xiaomimimo)
   381 points by thm 13 hours ago | hide | 143 comments

7. ▲ Google Play sees 47% decline in apps since start of last year (techcrunch.com)
   179 points by GeekyBear 3 hours ago | hide | 93 comments

8. ▲ Someone at YouTube needs glasses (jayd.ml)
   878 points by jaydenmilne 7 hours ago | hide | 494 comments

9. ▲ I created Perfect Wiki and reached $250k in annual revenue without investors (habr.com)
   524 points by sochix 14 hours ago | hide | 296 comments

10. ▲ Future of OSU Open Source Lab in Jeopardy (osuosl.org)
    72 points by aendruk 3 hours ago | hide | 18 comments

11. ▲ Show HN: Create your own finetuned AI model using Google Sheets (promptrepo.com)
    63 points by QueensGambit 6 hours ago | hide | 27 comments

12. ▲ DeepSeek-Prover-V2 (github.com/deepseek-ai)
    262 points by meetpateltech 5 hours ago | hide | 49 comments

# Martin Fowler Bliki

A website on building software effectively

https://martinfowler.com

Author

Works at ThoughtWorks

Exploring Generative AI

https://martinfowler.com/articles/exploring-gen-ai.html

The DeepSeek Series: A Technical Overview

https://martinfowler.com/articles/deepseek-papers.html

# ThoughtWorks Technology Radar

https://www.thoughtworks.com/radar

Volume 32 | April 2025

Techniques

Tools

Platforms

Languages & Frameworks

Adopt

Trial
    Worth pursing
    Try on projects that can handle risk

Assess
    Worth exploring
    How will it affect your enterprise

Hold
    Proceed with caution

## Techniques

**Adopt**
1. Data product thinking
2. Fuzz testing
3. Software Bill of Materials
4. Threat modeling

**Trial**
5. API request collection as API product artifact
6. Architecture advice process
7. GraphRAG
8. Just-in-time privileged access management
9. Model distillation
10. Prompt engineering
11. Small language models
12. Using GenAI to understand legacy codebases

**Assess**
13. AI-friendly code design
14. AI-powered UI testing
15. Competence envelope as a model for understanding system failures
16. Structured output from LLMs

**Hold**
17. AI-accelerated shadow IT
18. Complacency with AI-generated code
19. Local coding assistants
20. Replacing pair programming with AI
21. Reverse ETL
22. SAFe™

## Platforms

**Adopt**
23. GitLab CI/CD
24. Trino

**Trial**
25. ABsmartly
26. Dapr
27. Grafana Alloy
28. Grafana Loki
29. Grafana Tempo
30. Railway
31. Unblocked
32. Weights & Biases

**Assess**
33. Arize Phoenix
34. Chainloop
35. Deepseek R1
36. Deno
37. Graphiti
38. Helicone
39. Humanloop
40. Model Context Protocol (MCP)
41. Open WebUI
42. pg_mooncake
43. Reasoning models
44. Restate
45. Supabase
46. Synthesized
47. Tonic.ai
48. turbopuffer
49. VectorChord

**Hold**
50. Tyk hybrid API management

## Tools

**Adopt**
51. Renovate
52. uv
53. Vite

**Trial**
54. Claude Sonnet
55. Cline
56. Cursor
57. D2
58. Databricks Delta Live Tables
59. JSON Crack
60. MailSlurp
61. Metabase
62. NeMo Guardrails
63. Nyx
64. OpenRewrite
65. Plerion
66. Software engineering agents
67. Tuple
68. Turborepo

**Assess**
69. AnythingLLM
70. Gemma Scope
71. Hurl
72. Jujutsu
73. kubenetmon
74. Mergiraf
75. ModernBERT
76. OpenRouter
77. Redactive
78. System Initiative
79. TabPFN
80. v0
81. Windsurf
82. YOLO

## Languages and Frameworks

**Adopt**
83. OpenTelemetry
84. React Hook Form

**Trial**
85. Effect
86. Hasura GraphQL engine
87. LangGraph
88. MarkItDown
89. Module Federation
90. Prisma ORM

**Assess**
91. .NET Aspire
92. Android XR SDK
93. Browser Use
94. CrewAI
95. ElysiaJs
96. FastGraphRAG
97. Gleam
98. GoFr
99. Java post-quantum cryptography
100. Presidio
101. PydanticAI
102. Swift for resource-constrained applications
103. Tamagui
104. torchtune

**Hold**
105. Node overload

# The End (Almost)

Hope you learned a lot and found this course useful