## Al Cloud for Agent Development

3-dimensional optimization

Lin Qiao | CEO & Cofounder, Fireworks Al



## 2025 is the year of Agents

Coding Agents

Document Agents Sales/Marketing Agents

Hiring Agents

Customer Service Agents

Retail

Insurance

Finance

Education

Health/Medical

Manufacturing

Security



#### Fireworks AI



#### Fireworks AI



#### Fireworks AI



# cursor Uber Upath



## Principle 1:

#### **Model is Your Product**

with product-model codesign

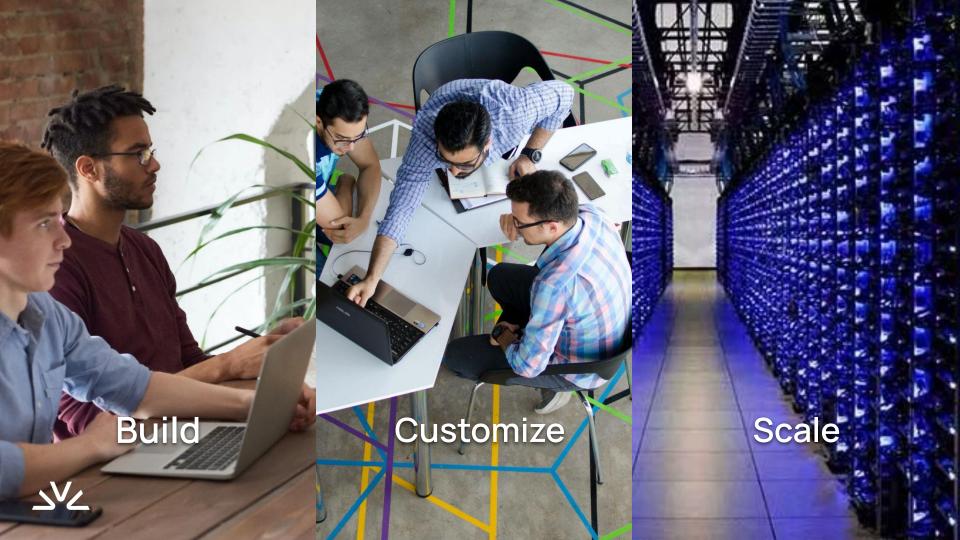


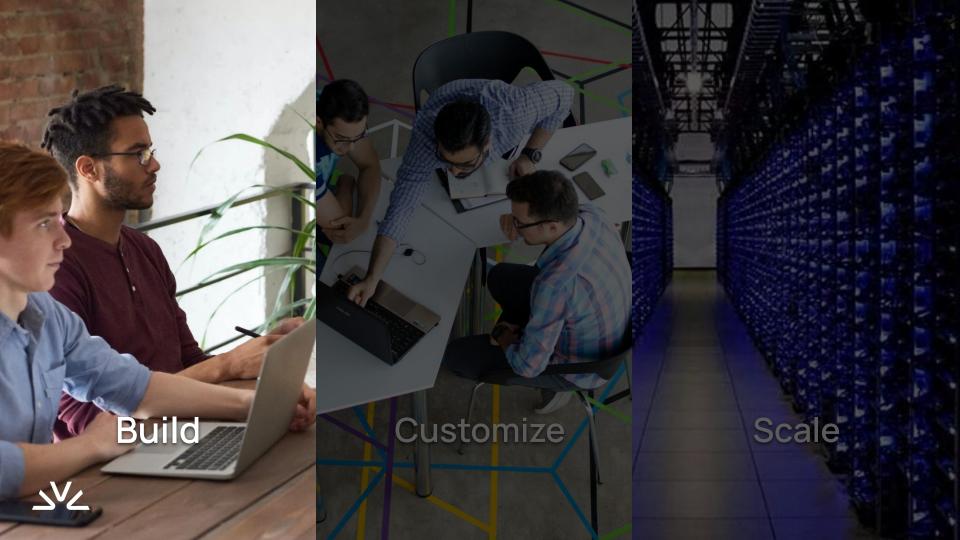
## Principle 2:

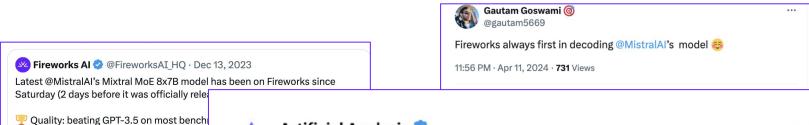
#### Model is Your IP

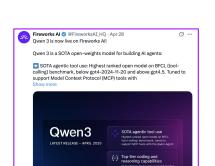
through data flywheel











Show more

Speed: fastest inference engine reaching



It is bizarre that third-party providers are hosting Gemma 2 before Google is.

We are now seeing third-party API providers host Gemma 2 with @FireworksAl\_HQ launching their Gemma 2 9B API.



t3. Fireworks Al



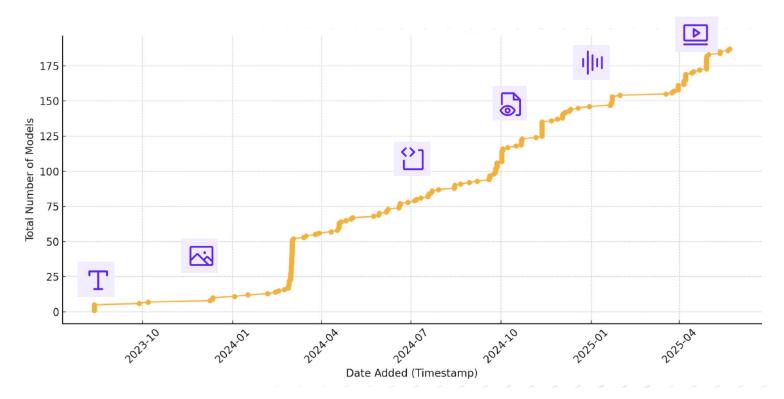


...



√ Fireworks AI

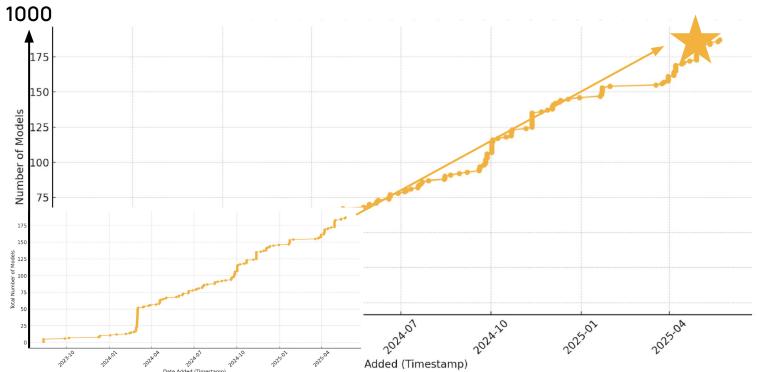
#### **SOTA Open Models Enabled**







#### Access to 1000+ models today!

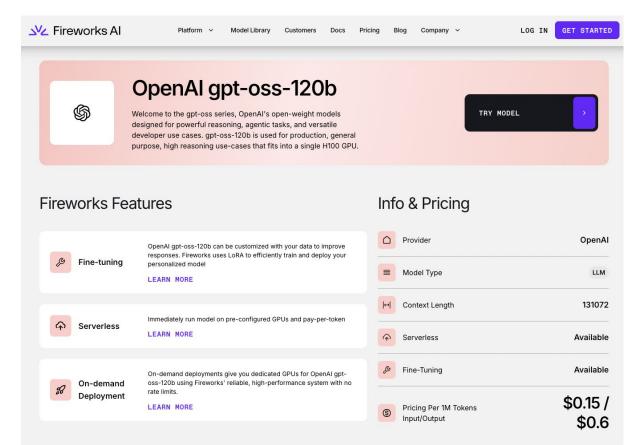




#### What a summer of open models!

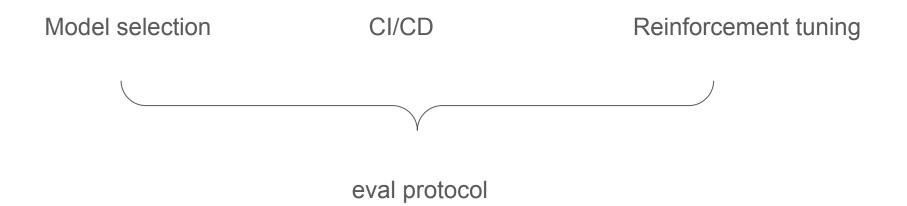
Model	Total Params, B	Active Params, B	KV cache entry, KB	Max Context	Layer s	Hidden Dim	Experts (MoE)			
							Total	Active	Spasity	Routing
gpt-oss-120b	120	5.1	36	131,072	36	2880	128	4	32	Simple
gpt-oss-20b	20	3.6	24	131,072	24	2880	32	4	8	Simple
DeepSeek V3/R1	671	37	69	163,840	61	7168	256	8	32	Grouped
Kimi K2	1000	32	69	131,072	61	7168	384	8	48	Simple
Qwen3-235B	235	22	188	262,144	94	4096	128	8	16	Simple
Qwen3-30B	30.5	3.3	96	262,144	48	2048	128	8	16	Simple
GLM-4.5	355	32	368	131,072	92	5120	160	8	20	Simple
GLM-4.5-Air	106	12	184	131,072	46	4096	128	8	16	Simple

## **Experiment Platform (GA)**



- SOTA Open models
- Fast, Powerful Tuning with Zero Setup
- Flexible Capacity & 1-click Deploy
- Private and Secure
- Build SDK experiment as code

## Evaluation is a big pain



#### <u>evalprotocol.io</u> github.com/eval-protocol



The open-source toolkit for building your internal model leaderboard.



Documentation

#### **What is Eval Protocol?**

When you have multiple Al models to choose from—different versions, providers, or configurations—how do you know whose for your use case?

#### Torganization Repositories

#### **Core Projects**

Repository	Description					
eval-protocol	Main specification, documentation, and examples					
python-sdk	Python implementation					

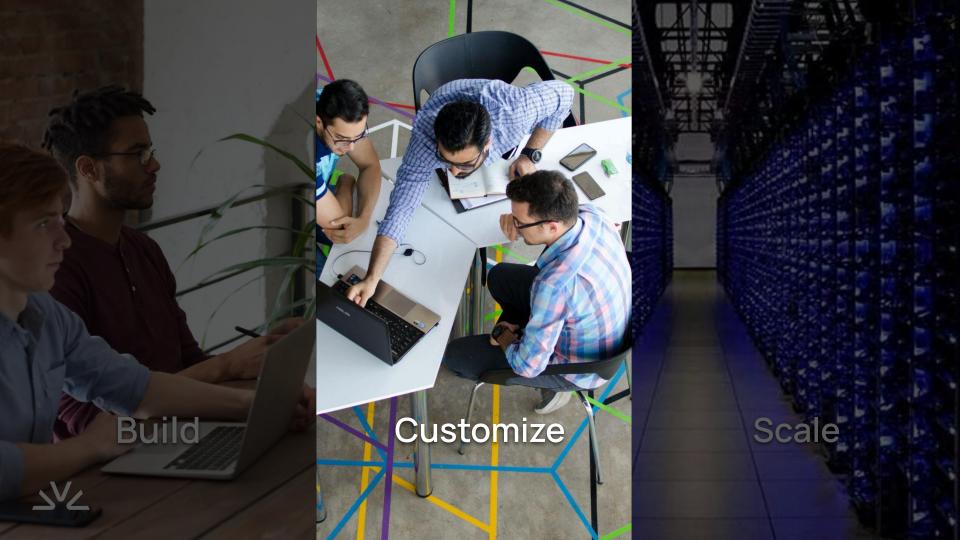
## Build with us

100X
Faster iteration

10K+
Companies

18X Growth in a year





## Customize with data flywheel

**Product** Model



# Supervised Fine Tuning update

- Big SOTA models
- Longer context
- Better quality for quantization
- Faster training



#### Announcing

## Reinforcement Fine Tuning (Beta)

- RFT SOTA open models
- Composing SFT and RL
- Flexible evaluation
- Ease of use with reward-kit SDK and web IDE



#### Supervised Fine Tuning V2

#### **Better Quality**

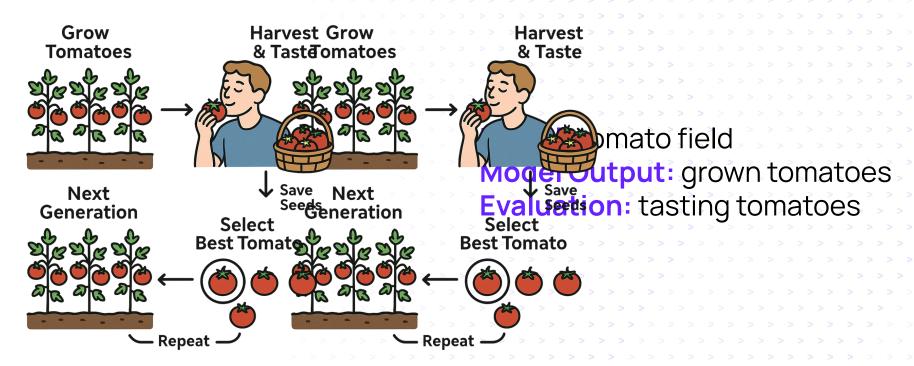
- Fine tune DeepSeek R1, V3, Llama and other SOTA models
- Access full 131K context window for training
- Preserve model quality with Quantization Aware Training

#### Faster Speed

- 2X faster training speed
- Better speed, quality with Multi-token Prediction Tuning



#### Reinforcement Fine Tuning





#### Reinforcement Fine Tuning Workflow

```
...
                                 Evaluator
def evaluate(messages: list[dict], ground_truth: int, **kwargs) -> dict:
   answer = messages[-1]['content']
    last_line = answer.split('\n')[-1]
   mg = re.match(r"Answer: (\d+)", last_line)
   score = 0.0
   if mg is not None:
        extracted_answer = int(mg.group(1))
        score = 1.0 if extracted_answer == ground_truth else 0.0
   return {
        "score": score,
        "reason": str(ground_truth),
        "is_score_valid": True,
```

Create Evaluator (Reward Function)

Select Model & Dataset

Run Training Job

Review Results and Deploy

## Reinforcement Fine Tuning

DEMO



Ranking

#### Better results than proprietary models

Candidate selection

9X Faster than o4-mini 20%
Better accuracy than o4-mini



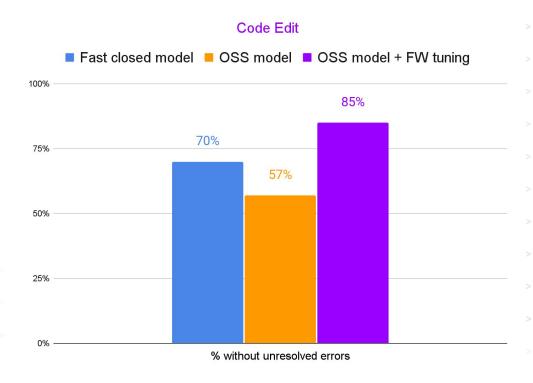
Code Edit

"The awesome news is our current model takes 2 passes to fix it, while using Fireworks this new model took 1. On a 800 LOC file, that is huge!"

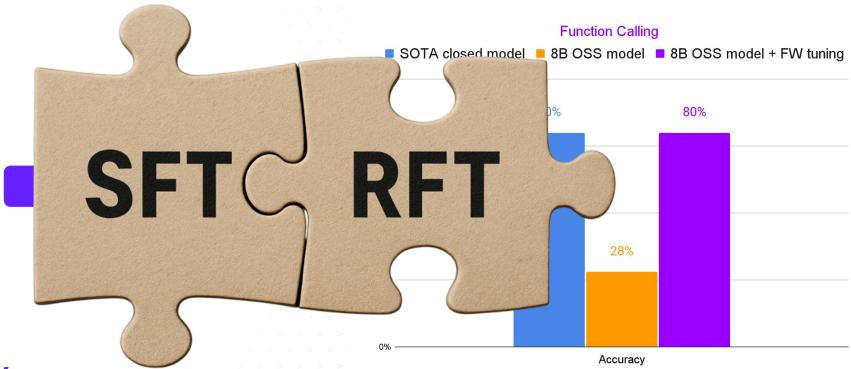
#### V

#### Better results than proprietary models

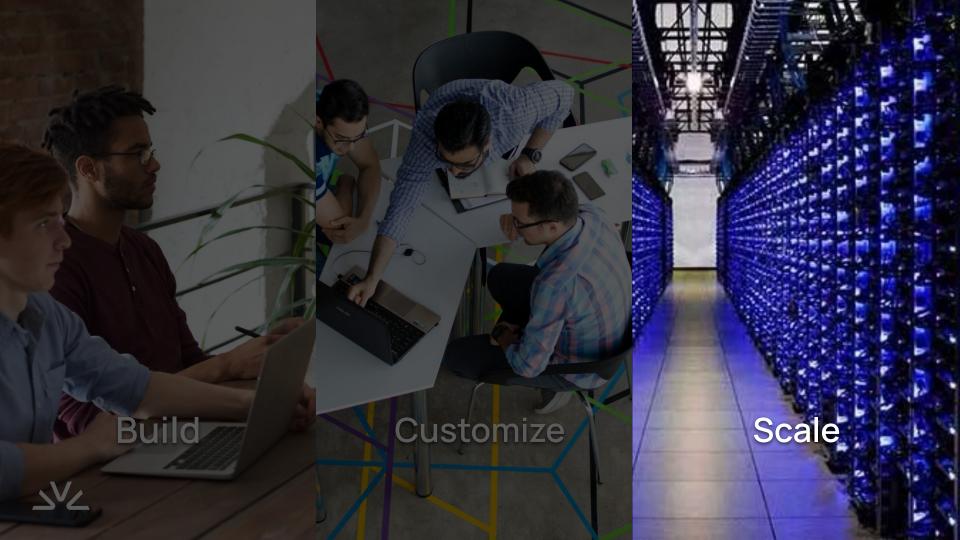
#### Code Rewrite for very large files



#### SOTA results with Fireworks Tuning







#### Scale Reliably with high performance

- ? Which hardware?
- ? Where is capacity?
- System failures with a thousand cuts
- ? How to optimize for production



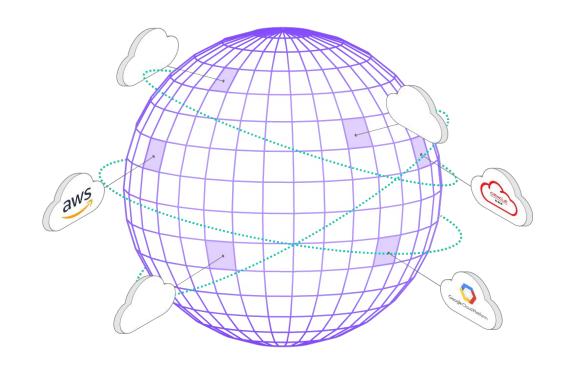
#### Global distribution

8

clouds

18

regions



+ bring your own cloud



#### Announcing

### Virtual Cloud

- Across 8 CSPs and SOTA hardware
- Disaggregated engine for max efficiency
- High reliability and global scheduling
- Enterprise grade privacy and compliance
- Bring your own cloud or bucket



#### 3-D Optimizer v2 for Quality, Speed & Concurrency

- 5 speculative decoding modes
- 8 quantization modes
- 7 hardware SKUs
- 4 sharding schemes
- 5 multi-host setups
- 10+ kernel options
- 4 quality tuning approaches

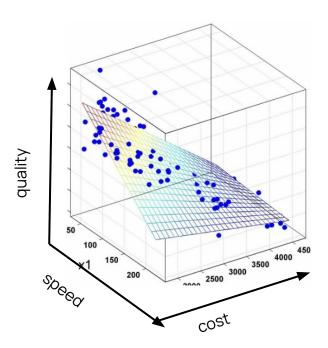
100,000+

possible



#### 3-D Optimizer for Quality, Speed & Concurrency

All in One





# 10,000,000,000,000+

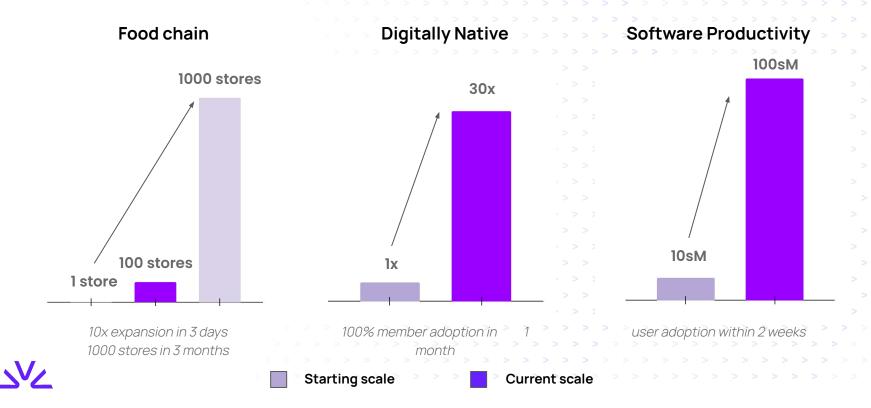
tokens / day

180,000+

requests / second



#### Enterprise customers scale fast



## Start building today!

- \$5k credits
- New product links

