# Alvin Wan

OpenAI research scientist optimizing inference for Large Language Models

## University of California, Berkeley

**Ph.D. in Computer Science** · [4k+ citations](#) · [1.8k+ stars](#) · NSF GRFP 2018[1]     **2018 - 2022**

**B.S. in Electrical Engineering and Computer Science** · 3.81 Major GPA     **2014 - 2018**
Samuel Silver Memorial Award ('18), Summer Undergraduate Research Fellowship ('17),
Leadership Award ('16), Dean's Honor List ('16), Regents' and Chancellor's Scholarship ('14)

**Lecturer** 2x, **Head Student Instructor** 6x · 4.83 / 5.00 rating     **2015 - 2021**
Taught Machine Learning, Discrete Mathematics, Data Science. Managed 70+ person staffs
serving 500-1100 students. Wrote [math](#), [probability](#), [compsci](#) booklets (10k+ downloads)

## OpenAI

**Member of Technical Staff**     Oct '24 - **now**
Deployed multimodal inference optimizations for OpenAI's largest model to date, [GPT 4.5](#).
Enabled doubled context length for [GPT 4.1](#).

## Apple

**Senior Research Scientist**     May '22 - Oct '24
Trained and deployed [Apple Intelligence](#)'s first low-precision foundation models. Deployed
M*-specific optimizations. Facilitated internships, UC Berkeley sponsorship. [ICML'23](#)

## Tesla *AutoPilot*

**AI Intern**     Mar - Aug '21
Trained and deployed two of Tesla's first FSD radar-less perception models to millions of
cars, improving safety-critical kinematics predictions for pedestrians, cyclists etc.

## Facebook *Reality Labs*

**Research Intern** - Neural Architecture Search, [FBNetV2](#), [FBNetV3](#) @ CVPR     May '19 - Mar '21
**Software Engineering Intern** - Halved media requests for Android messenger     May - Aug '16

## Startups

- **Full-Time** Machine Learning Engineer at REX Homes     May '18 - Aug '21
- Machine Learning Intern at DeepScale *(acquired by Tesla)*     May - Aug '17
- Software Engineering Intern at Getexp     May - Aug '15

## International Awards

- Microsoft Imagine Cup "Big Data" **World Finals Top 6 Finalists** (2018, of 40k+ entrants, 200+ countries)
- Rookies Co. "Web and Mobile" **Int'l Top 16 Semifinalist** (2017, ~9k entries, 80+ countries)
- Adobe Design Achievement Awards "Social Impact" **Int'l Semifinalist** (2017, 2018, ~7k entrants)

---

[1] One of 20 nation-wide recipients for Machine Learning, representing ~0.1% of 12,000+ applicants.