# JACOB KAHN

## PERSONAL INFORMATION

| | |
|---:|:---|
| *Email* | jacobkahn1@gmail.com |
| *Website* | jacobkahn.me |
| *GitHub* | github.com/jacobkahn |
| *Scholar* | Google Scholar |

*Research Manager and Engineer at FAIR, Meta AI, developing scalable systems and models for reasoning, code generation, and multimodal learning.*

RESEARCH AREAS: Deep learning, distributed systems, reasoning, code generation, multimodal learning, and speech.

## EXPERIENCE

**2018 - present**

RESEARCH MANAGER AND ENGINEER, FAIR, META AI
*Menlo Park, California; New York City, New York.*
Technical leader across reasoning, code generation, and multimodal modeling, working full-stack on model–system co-design for scaling deep learning. Recruited and mentored researchers to senior and staff roles at Meta.

- REASONING AND CODE GENERATION: Co-leading Code World Model (CWM). Direction for scaling, data generation, RL infrastructure, and inference; leading a team of scientists and engineers; steering over 30 ExaFLOPs of compute.
- MULTIMODAL MODELING: Led inference, infrastructure, and scaling for Chameleon and contributed to Transfusion. Drove projects in end-to-end speech recognition, including Libri-Light and wav2letter.
- OPEN-SOURCE FRAMEWORKS: Led development of Flashlight, co-led Shumai and wav2letter, projects that influenced PyTorch `distributed`, `compile`, and Meta's ML infrastructure roadmaps.
- HARDWARE AND SCALING STRATEGY: Planning for Meta's GPU clusters and shaping hardware strategy for large-model scaling; guiding multi-billion-dollar compute investments.

**2024 - present**

COMPUTER SCIENCE FACULTY, UNIVERSITY OF PENNSYLVANIA
*Philadelphia, Pennsylvania.*
Affiliated with Penn's NetDB Laboratory.
Design and teach *CIS 5690: GPU Programming and Machine Learning Systems*.

**2016 - 2018**

Engineering Intern, FACEBOOK
*Menlo Park, California.*
Built large, globally distributed systems and algorithms for high-performance stream processing in C++.

## EDUCATION

**2016 - 2018**

M.S.E. IN COMPUTER AND INFORMATION SCIENCE, University of Pennsylvania.

**2014 - 2018**

JEROME FISHER PROGRAM IN MANAGEMENT AND TECHNOLOGY
University of Pennsylvania

B.S.E. IN COMPUTER AND INFORMATION SCIENCE, Penn Engineering
THESIS: *Computer Vision for Multiplayer Anchoring in Real-Time Augmented Reality Systems*

B.S. IN ECONOMICS with concentrations in OPERATIONS RESEARCH and MANAGEMENT, The Wharton School

*Full list available on Google Scholar.*

2025
CWM: An Open-Weights LLM for Research on Code Generation with World Models
CODE WORLD MODEL TEAM, **(Core Team, Frameworks & Infrastructure Lead)**
*arXiv Preprint 2510.02387*

Efficient Hardware Scaling and Diminishing Returns in Large-Scale Training of Language Models
Jared FERNANDEZ, Luca WEHRSTEDT, Leonid SHAMIS, Mostafa ELHOUSHI, Kalyan SALADI, Yonatan BISK, Emma STRUBELL, **Jacob Kahn**
In *Transactions on Machine Learning Research, 2025*

Transfusion: Predict the Next Token and Diffuse Images with One Multi-Modal Model
Chunting ZHOU, Lili YU, Arun Babu, Kushal TIRUMALA, Michihiro YASUNAGA, Leonid SHAMIS, **Jacob Kahn**, Xuezhe MA, Luke ZETTLEMOYER, Omer LEVY
In *International Conference on Learning Representations (Oral)* Singapore, 2025

2024
Chameleon: Mixed-modal Early-Fusion Foundation Models
CHAMELEON TEAM **(Evaluation & Inference Lead)**
*arXiv Preprint 2405.09818*

2023
Reasoning over Public and Private Data in Retrieval-Based Systems
Simran ARORA, Patrick LEWIS, Angela FAN, **Jacob Kahn\***, Christopher RÉ\*
in the *Transactions of the Association for Computational Linguistics (TACL), presented at the Proceedings of the ACL, Toronto, Canada, 2023.*
also in *AAAI — workshop on Knowledge Augmented Methods for NLP (Oral), Washington D.C. 2023*
\* = Equal advising.

RA-DIT: Retrieval-Augmented Dual Instruction Tuning
X. LIN, X. CHEN, M. CHEN, W. SHI, M. LOMELI, R. JAMES, P. RODRIGUEZ, **J. Kahn**, G. SZILVASU, M. LEWIS, L. ZETTLEMOYER, S. YIH
In *International Conference on Learning Representations (ICLR), Vienna, Austria, 2024*

2022
Flashlight: Enabling Innovation in Tools for Machine Learning
**J. Kahn**, V. PRATAP, T. LIKHOMANENKO, Q. XU, A. HANNUN, J. CAI, P. TOMASELLO, A. LEE, E. GRAVE, G. AVIDOV, B. STEINER, V. LIPTCHINSKY, G. SYNNAEVE, R. COLLOBERT
In *International Conference on Machine Learning (ICML) (Spotlight), Baltimore, Maryland, 2022*

2021
slimIPL: Language-Model-Free Iterative Pseudo-Labeling
Tatiana LIKHOMANENKO\*, Qiantong XU\*, **Jacob Kahn**, Gabriel SYNNAEVE, Ronan COLLOBERT
In *Proceedings of Interspeech, Brno, Czech Republic, 2021*
\* = Equal contribution.

2020
Self-Training for End-to-End Speech Recognition
**Jacob Kahn**, Ann LEE, Awni HANNUN
In *Proc. of the 45th IEEE International Conference in Acoustic, Speech and Signal Processing (ICASSP), Barcelona, Spain, 2020.*

Libri-Light: A Benchmark for ASR with Limited or No Supervision
**J. Kahn**\*, M. RIVIÈRE\*, W. ZHENG\*, E. KHARITONOV\*, Q. XU\*, P.E. MAZARÉ\*, J. KARADAYI\*, V. LIPTCHINSKY, R. COLLOBERT, C. FUEGEN, T. LIKHOMANENKO, G. SYNNAEVE, A. JOULIN, A. MOHAMED, E. DUPOUX
In *Proc. of the 45th IEEE International Conference in Acoustic, Speech and Signal Processing (ICASSP), Barcelona, Spain, 2020.*
\* = Equal contribution.

End-to-End ASR: from Supervised to Semi-Supervised Learning with Modern Architectures
G. SYNNAEVE\*, Q. XU\*, **J. Kahn**\*, E. GRAVE\*, T. LIKHOMANENKO\*, V. PRATAP, A.

Sriram, V. Liptchinsky, R. Collobert
In *ICML Workshop on Self-Supervision in Audio and Speech, Virtual, 2020*
* = Equal contribution.

Differentiable Weighted Finite-State Transducers
Awni Hannun, Vineel Pratap, **Jacob Kahn**, Wei-Ning Hsu
*arXiv Preprint 2010.01003*

2019     wav2letter++: A Fast Open-source Speech Recognition System
Vineel Pratap, Awni Hannun, Qiantong Xu, Jeff Cai, **Jacob Kahn**, Gabriel Synnaeve, Vitaliy Liptchinsky, Ronan Collobert
In *The 44th IEEE International Conference in Acoustic, Speech and Signal Processing (ICASSP), Brighton, UK, 2019.*

## INVITED TALKS

2025     *Tracing Program Execution for Neural Planning and Reasoning*
AI Engineer, New York, NY, scheduled for November 2025

*Scaling Asynchronous RL Recipes for Learning to Code*
Meta Superintelligence Labs, St. Louis, MO, scheduled for November 2025

2024     *Trends in Deep Learning Computation and Frameworks*
Intel, Santa Clara, CA, September 2024

*Multimodal Generative Models: Chameleon and Beyond*
Charles River Ventures and Greylock, San Francisco, CA, August 2024

*GPU Computing at Scale*
Guest Lecture, CIS 565: GPU Programming & Architecture, Penn, Philadelphia, PA, July 2024

2023     *Shumai: Fast, Flexible Machine Learning in TypeScript*
Google DeepMind, Mountain View, CA, April 2023

2020     *Scaling Deep Learning for Automatic Speech Recognition*
NVIDIA GPU Technology Conference, San Jose, CA, March 2020

## TEACHING AND SERVICE

*University of Pennsylvania*

GPU Computing for Machine Learning Systems · CIS 5690 · Instructor
Created and teach a graduate course on GPU programming with CUDA, large-scale deep learning systems, performance engineering, and compilers. Project-based, in CUDA and C++.

Head TA (as student): Algorithms (CIS 320), Graduate Operating Systems (CIS 548)
TA Experience: Undergraduate and graduate courses in distributed systems, computer architecture, and databases, all in C/C++.

*Reviewing*     Regular reviewer at NeurIPS, ICLR, and ICML.