

# Dan Fu

---

## Contact

E-mail: [danfu@cs.stanford.edu](mailto:danfu@cs.stanford.edu)  
Homepage: <https://www.danfu.org/>

## Research Interests

My research interests are in the intersection between machine learning and systems. I'm particularly interested in developing more efficient algorithms and architectures for ML, along with hardware-aware systems solutions to make them practically effective.

## Education

2018–	PhD in Computer Science, Stanford University Advisors: Christopher Ré, Kayvon Fatahalian
2014–2018	AB/SM in Computer Science, Harvard University Thesis: <i>Design of Influencing Agents for Flocking in Low-Density Settings</i>

## Experience

2018–	PhD Student, <b>Stanford University</b>
2022–	Academic Partner, <b>Together AI</b>
2020–	Founder/Podcast Host, <b>Stanford MLSys Seminar</b> 16,000+ subscribers, 30,000+ monthly views
Summer 2019	Research Intern, <b>Argo AI</b>
Summer 2016, 2017	Software Engineering Intern, <b>Google</b>
Summer 2015	Field Engineering Intern, <b>Tamr</b>
Summer 2014	Software Engineering Intern, <b>Interactive Intelligence</b>
Summer 2013	Development Intern, <b>DyKnow</b>

## Awards

May 2024	<b>Stanford Data Science Open Source Software Prize 2024</b> Inaugural open source software prize awarded for FlashAttention
December 2023	<b>Best Poster: Efficient Natural Language and Speech Processing Workshop at NeurIPS 2023</b> FlashFFTConv: Efficient Convolutions for Long Sequences with Tensor Cores
December 2023	<b>Oral Presentation: NeurIPS 2023</b> Monarch Mixer: A Simple Sub-Quadratic GEMM-Based Architecture
July 2023	<b>Oral Presentation: ICML 2023</b> Hyena Hierarchy: Towards Larger Convolutional Language Models High-throughput Generative Inference of Large Language Models with a Single GPU
May 2023	<b>Top-25% / Spotlight: ICLR 2023</b> Hungry Hungry Hippos: Towards Language Modeling with State Space Models
August 2022	<b>Best Student Paper Runner Up: UAI 2022</b> Shoring Up the Foundations: Fusing Model Embeddings and Weak Supervision
July 2022	<b>NDSEG Award for Exemplary Impact and Relevance to DoD Research Objectives</b> Awarded for research presented at 2022 NDSEG conference
July 2022	<b>Best Paper: Hardware Aware Efficient Training Workshop at ICML 2022</b> FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness
February 2022	<b>Best Paper: AIBSD Workshop at AAAI 2022</b>

2019-2022	The Details Matter: Preventing Class Collapse in Supervised Contrastive Learning <b>United States Department of Defense NDSEG Fellow</b>
2019-2020	<b>Brown Institute for Media Innovation Magic Grant</b>
Fall 2017, Spring 2018	<b>Harvard Derek Bok Center Certificate of Distinction in Teaching</b>
2014	<b>Presidential Scholar</b>
2012	<b>Siemens Research Competition National Runner-Up</b>

## Publications

- [1] **Benchmarking and Building Long-Context Retrieval Models with LoCo and M2-BERT**  
*International Conference on Machine Learning (ICML) (2024) and ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models (ME-FoMo) (2024)*  
Jon Saad-Falcon, **Daniel Y. Fu**, Simran Arora, Neel Guha, Christopher Ré
- [2] **FlashFFTConv: Efficient Convolutions for Long Sequences with Tensor Cores**  
*International Conference on Learning Representations (ICLR) (2024)*  
*Workshop on Efficient Natural Language and Speed Processing at NeurIPS (2023), **Best Poster***  
**Daniel Y. Fu\***, Hermann Kumbong\*, Eric Nguyen, and Christopher Ré
- [3] **Monarch Mixer: A Simple Sub-Quadratic GEMM-Based Architecture**  
*Advances in Neural Information Processing Systems (NeurIPS) (2023), **Oral Presentation***  
**Daniel Y. Fu**, Simran Arora, Jessica Grogan, Isys Johnson, Sabri Eyuboglu, Armin W. Thomas, Benjamin F. Spector, Michael Poli, Atri Rudra, and Christopher Ré
- [4] **Laughing Hyena Distillery: Extracting Compact Recurrences from Convolutions**  
*Advances in Neural Information Processing Systems (NeurIPS) (2023)*  
Stefano Massaroli\*, Michael Poli\*, **Daniel Y. Fu\***, Hermann Kumbong, David W. Romero, Rom Nishijima Parnichkun, Aman Timalsina, Quinn McIntyre, Beidi Chen, Atri Rudra, Ce Zhang, Christopher Ré, Stefano Ermon, and Yoshua Bengio
- [5] **High-throughput Generative Inference of Large Language Models with a Single GPU**  
*International Conference on Machine Learning (ICML) (2023), **Oral Presentation***  
Ying Sheng, Lianmin Zheng, Binhang Yuan, Zhuohan Li, Max Ryabinin, **Daniel Y. Fu**, Zhiqiang Xie, Beidi Chen, Clark Barrett, Joseph E. Gonzalez, Percy Liang, Christopher Ré, Ion Stoica, Ce Zhang
- [6] **Hyena Hierarchy: Towards Larger Convolutional Language Models**  
*International Conference on Machine Learning (ICML) (2023), **Oral Presentation***  
Michael Poli, Stefano Massaroli, Eric Nguyen, **Daniel Y. Fu**, Tri Dao, Stephen Baccus, Yoshua Bengio, Stefano Ermon, Christopher Ré
- [7] **Simple Hardware-Efficient Long Convolutions for Sequence Modeling**  
*International Conference on Machine Learning (ICML) and ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models (ME-FoMo) (2023)*  
**Daniel Y. Fu\***, Elliot L. Epstein\*, Eric Nguyen, Armin W. Thomas, Michael Zhang, Tri Dao, Atri Rudra, Christopher Ré
- [8] **Hungry Hungry Hippos: Towards Language Modeling with State Space Models**  
*International Conference on Learning Representations (ICLR) (2023), **Notable Top-25% / Spotlight Presentation***  
**Daniel Y. Fu\***, Tri Dao\*, Khaled K. Saab, Armin W. Thomas, Atri Rudra, Christopher Ré
- [9] **FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness**  
*Neural Information Processing Systems (NeurIPS) (2022)*  
*Sparsity in Neural Networks: Advancing Understanding and Practice (2022), **Oral Presentation***  
*Hardware Aware Efficient Training Workshop at ICML (2022), **Best Paper***  
Tri Dao, **Daniel Y. Fu**, Stefano Ermon, Atri Ruda, Christopher Ré
- [10] **Shoring Up the Foundations: Fusing Model Embeddings and Weak Supervision**  
*Conference on Uncertainty in Artificial Intelligence (UAI) (2022), **Best Student Paper Runner Up***  
Mayee F. Chen\*, **Daniel Y. Fu\***, Dyah Adila, Michael Zhang, Frederic Sala, Kayvon Fatahalian, Christopher Ré

- [11] **Perfectly Balanced: Improving Transfer and Robustness in Supervised Contrastive Learning**  
*International Conference on Machine Learning (ICML)* (2022)  
 Workshop version: **The Details Matter: Preventing Class Collapse in Supervised Contrastive Learning**  
*Workshop on Artificial Intelligence with Biased or Scarce Data (AIBSD) at AAAI* (2022), **Best Paper**  
 Mayee F. Chen\*, **Daniel Y. Fu\***, Avanika Narayan, Michael Zhang, Zhao Song, Kayvon Fatahalian, Christopher Ré
- [12] **TABi: Type-Aware Bi-encoders for End-to-End Entity Retrieval**  
*Findings of the Association for Computational Linguistics: ACL* (2022)  
 Megan Leszczynski, **Daniel Y. Fu**, Mayee F. Chen, Christopher Ré
- [13] **Harmonizing Attention: Attention Map Consistency For Unsupervised Fine-Tuning**  
*Bridging the Gap: From Machine Learning Research to Clinical Practice Workshop at NeurIPS* (2021)  
 Ali Mirzazadeh, Florian Dubost, Maxwell Pike, Krish Maniar, **Daniel Y. Fu**, Khaled K Saab, Christopher Lee-Messer, Daniel Rubin
- [14] **Analyzing Who and What Appears in a Decade of US Cable TV News**  
*ACM SigKDD Conference on Knowledge Discovery & Data Mining (KDD)* (2021)  
 James Hong, Will Crichton, Haotian Zhang, **Daniel Y. Fu**, Jacob Ritchie, Jeremy Barenholtz, Ben Hannel, Xinwei Yao, Michaela Murray, Geraldine Moriba, Maneesh Agrawala, Kayvon Fatahalian
- [15] **Beyond the Pixels: Exploring the Effect of Video File Corruptions on Model Robustness**  
*ECCV 2020 Workshop on Adversarial Robustness in the Real World* (2020)  
 Trenton Chang, **Daniel Y. Fu**, Sharon Yixuan Li, Christopher Ré
- [16] **Fast and Three-rious: Speeding Up Weak Supervision with Triplet Methods**  
*International Conference on Machine Learning (ICML)* (2020)  
**Daniel Y. Fu\***, Mayee F. Chen\*, Frederic Sala, Sarah M. Hooper, Kayvon Fatahalian, Christopher Ré
- [17] **Multi-Resolution Weak Supervision for Sequential Data**  
*Neural Information Processing Systems (NeurIPS)* (2019)  
 Frederic Sala, Paroma Varma, Jason Fries, **Daniel Y. Fu**, Shiori Sagawa, Saelig Khattar, Ashwini Ramamoorthy, Ke Xiao, Kayvon Fatahalian, James R. Priest, Christopher Ré
- [18] **Video Event Specification Using Programmatic Composition**  
*AI Systems Workshop at SOSIP* (2019), **Oral Presentation**  
**Daniel Y. Fu**, Will Crichton, James Hong, Xinwei Yao, Haotian Zhang, Anh Truong, Avanika Narayan, Maneesh Agrawala, Christopher Ré, Kayvon Fatahalian
- [19] **Influencing Flock Formation in Low-Density Settings**  
*International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)* (2018)  
**Daniel Y. Fu**, Emily S. Wang, Peter M. Krafft, Barbara J. Grosz

## Journal Publications

- [20] **Orexinergic Neurotransmission in Temperature Responses to Methamphetamine and Stress: Mathematical Modeling as a Data Assimilation Approach**  
*PLOS ONE*, May 20, 2015  
 Abolhassan Behrouzvaziri, **Daniel Y. Fu**, Patrick Tan, Yeonjoo Yoo, Maria V. Zaretskaia, Daniel E. Rusyniak, Yaroslav I. Molkov, Dmitry V. Zaretsky
- [21] **Chaos and Robustness in a Single Family of Genetic Oscillatory Networks**  
*PLOS ONE*, March 25, 2014  
**Daniel Y. Fu**, Patrick Tan, Alexey Kuznetsov, Yaroslav I. Molkov

## Unpublished Work and Preprints

- [22] **Rekall: Specifying Video Events using Compositions of Spatiotemporal Labels**  
*arXiv preprint arXiv:1910.02993*, October 2019  
**Daniel Y. Fu**, Will Crichton, James Hong, Xinwei Yao, Haotian Zhang, Anh Truong, Avanika Narayan, Maneesh Agrawala, Christopher Ré, Kayvon Fatahalian

[23] **Automatic Parallelization of Sequential Programs**

*arXiv preprint arXiv:1809.07684*, July 2018

Peter Kraft, Amos Waterland, **Daniel Y. Fu**, Anitha Gollamudi, Shai Szulanski, Margo Seltzer

**Invited Talks**

- 2024 **The Unreasonable Power of Synthetics for Efficient Machine Learning**  
Young Professional Symposium at MLSys 2024, Santa Clara, CA
- 2024 **Hardware-Aware Efficient Primitives for Machine Learning**  
Duke University, Durham, NC  
University of Michigan, Ann Arbor, MI  
Cornell Tech, New York, NY  
Yale University, New Haven, CT  
UCLA, Los Angeles, CA  
NYU New York, NY  
UT Austin, Austin, TX  
Harvard University, Cambridge, MA  
Purdue University, West Lafayette, IN  
UCSD, San Diego, CA  
Cornell University, Ithaca, NY  
Northwestern University, Evanston, IL  
University of Wisconsin-Madison IDEA Seminar, Madison, WI  
Google, Mountain View, CA  
Apple, Virtual
- 2023 **Efficient Sub-Quadratic Architectures for Machine Learning**  
Bangkok AI Hack 2023, Thailand
- 2023 **Monarch Mixer: A Simple Sub-Quadratic GEMM-Based Architecture**  
SystemX Conference, Stanford, CA
- 2023 **Hungry Hungry Hippos**  
Neural Notes Podcast, Virtual  
AI Pub Deep Papers Podcast, Virtual  
IBM Research, Virtual
- 2023 **Perspectives in Generative AI Panel**  
Pear VC Firm, Menlo Park, CA
- 2022-2023 **FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness**  
Google, Mountain View, CA  
Stanford CS 217 Guest Lecture, Stanford, CA  
MLOps World: Machine Learning in Production, Virtual  
Meta PyTorch Performance Team, Virtual  
MosaicML, Virtual  
Google N2Formal Team, Mountain View, CA
- 2022 **Improving Transfer and Robustness of Supervised Contrastive Learning**  
Stanford MLSys Seminar, Stanford, CA  
KitWare Vision Research Group, Virtual
- 2019 **Rekall: Modeling Concepts in Video with Compositions of Spatiotemporal Labels**  
Intel, Bend, OR  
Stanford Graphics Group, Stanford, CA  
Stanford Vision Lab, Stanford, CA  
Stanford DAWN Retreat, Menlo Park, CA

## Open-Source Artifacts

2023	<b>FlashFFTConv: Efficient Convolutions for Long Sequences</b> A library for fast exact convolutions optimized for GPU. Now being used in production by Together AI to train and serve long-sequence convolutional language models. <a href="https://github.com/HazyResearch/flash-fft-conv">https://github.com/HazyResearch/flash-fft-conv</a>
2023	<b>Monarch Mixer BERT Models</b> A suite of BERT models trained with Monarch Mixer, from 80M to 341M parameters, supporting sequence lengths up to 8K. Now being served by Together AI and in use by MongoDB. <a href="https://github.com/HazyResearch/m2">https://github.com/HazyResearch/m2</a>
2023	<b>RedPajama-1T</b> A trillion+ token dataset for training large language models, mimicking the data gathering process from Llama-1. So far downloaded 1 million+ times. <a href="https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T">https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T</a>
2023	<b>Safari: Convolutions for Sequence Modeling</b> A training repository for gated convolution models on language, images, and long-sequence data. Used to train H3, Hyena, and M2. <a href="https://github.com/HazyResearch/safari">https://github.com/HazyResearch/safari</a>
2023	<b>Hungry Hungry Hippos (H3) Models</b> A suite of hybrid attention + gated SSM architectures trained on language modeling, up to 2.7B parameters. <a href="https://github.com/HazyResearch/H3">https://github.com/HazyResearch/H3</a>
2022	<b>FlashAttention</b> Fast and memory-efficient exact attention with IO-Awareness. Now integrated into PyTorch and used in every major AI research lab in industry. <a href="https://github.com/HazyResearch/flash-attention">https://github.com/HazyResearch/flash-attention</a>
2020	<b>FlyingSquid</b> A fast algorithm for weak supervision without SGD using method-of-moments estimation. In use at Snorkel AI. <a href="https://github.com/HazyResearch/flyingsquid">https://github.com/HazyResearch/flyingsquid</a>
2019	<b>Rekall: Compositional Video Event Specification</b> A library for analyzing video data using compositions of image labels. Once used at Argo AI for event mining. <a href="https://github.com/scanner-research/rekall">https://github.com/scanner-research/rekall</a>

## Teaching Experience

Fall 2021, Winter/Spring 2022, Fall 2023	<b>Instructor</b> , CS 528: Machine Learning Systems Seminar Stanford University
Winter 2023	<b>Interviewer</b> , CS 324: Advances in Foundation Models Stanford University
Spring 2018	<b>Teaching Fellow</b> , CS 152: Programming Languages Harvard University
Fall 2015, 2016, 2017	<b>Teaching Fellow</b> , CS 61: Systems Programming and Machine Organization Harvard University

## Mentorship

2024	<b>Aaryan Singhal</b> (Stanford CS undergrad)
2023	<b>Jon Saad-Falcon</b> (Stanford CS PhD rotator)
2023	<b>Hermann Kumbong</b> (Stanford CS MS)
2023	<b>Pranav Vaid</b> (Stanford CS undergrad/coterm)
2022-2023	<b>Elliot Epstein</b> (Stanford ICME PhD rotator)
2020	<b>Heidi Chen</b> (Stanford CS undergrad, now Google)
2019-2020	<b>Trenton Chang</b> (Stanford undergrad/coterm in American Studies, now EECS PhD at University of Michigan)
2019-2022	<b>Avanika Narayan</b> (Stanford CS undergrad/coterm, now CS PhD at Stanford)
2020-2021	<b>Manasi Ganti</b> (high school student, now CS undergrad at University of Washington)

## Service

2020–	<b>Creator/Host: Stanford MLSys Seminar Series</b>
ICML 2023, 2024	<b>Workshop Organizer: Efficient Systems for Foundation Models (ES-FoMo)</b>
2020-2023	<b>Stanford PhD Admissions Committee</b>

Referee/program committee member for NeurIPS 2024, ICLR 2023 Blog Post Track, NeurIPS 2023, ICML 2023, ICLR 2023, NeurIPS 2022, ICML 2022 (**top 10% reviewer, session chair**), ICLR 2022, NeurIPS 2021, ICML 2021 Workshop on ML for Data, Workshop on Weakly Supervised Learning @ ICLR 2021.

Last updated: June 13, 2024<sup>\*</sup>

---

<sup>\*</sup>CV template inspired by [Neel Guha](#) and [Christopher Morris](#).