EMILY CHENG

emilyshana.cheng@upf.edu

EDUCATION

Universitat Pompeu Fabra

PhD Candidate in Linguistics Linguistic complexity and its representation in biological and artificial brains Supervisor: Marco Baroni

Massachusetts Institute of Technology Master of Engineering in Computer Science (2022)

Supervisors: Boris Katz and Andrei Barbu Bachelor of Science in Computer Science and Engineering (2020) Bachelor of Science in Mathematics (2020)

RESEARCH SUMMARY

I study how structure and parsimony arise in representations of language. Using tools from manifold learning (representational level) and information theory (behavioral level), I'm interested in the interaction between representational geometry and functional demands of language. I'm especially interested in how low-dimensional geometry allows biological/artificial brains to handle the vast complexity of language.

PUBLICATIONS

- 1. Emily Cheng, Chris Wang, Greta Tuckute, Marco Baroni, Andrei Barbu. Mapping linguistic complexity in the brain. *In prep.*
- Emily Cheng, Carmen Amo Alonso, Marco Baroni. Linearly Controlled Language Generation with Performative Guarantees. In prep. Workshop version in MINT@NeurIPS 2024.
- 3. Emily Cheng^{*} and Richard Antonello^{*}. Evidence from fMRI supports a two-phase abstraction process in language models. *In prep.* Workshop version in UniReps@NeurIPS 2024. Awarded best abstract. (Oral)
- 4. Emily Cheng and Francesca Franzon. Principles of semantic and functional efficiency in grammatical patterning. Preprint, 2024. Under final review at PNAS.
- Jin Hwa Lee*, Thomas Jiralerspong*, Jade Yu, Yoshua Bengio, Emily Cheng. Geometric signatures of compositionality across a language model's lifetime. In Proceedings of ACL 2025. (Oral) Workshop version in NeurReps@NeurIPS 2024.
- 6. Emily Cheng, Diego Doimo, Corentin Kervadec, Iuri Macocco, Jade Yu, Alessandro Laio, Marco Baroni. Emergence of a High-Dimensional Abstraction Phase in Language Transformers. In Proceedings of ICLR 2025. (Poster)
- 7. Emily Cheng, Corentin Kervadec, Marco Baroni. Bridging Information-Theoretic and Geometric Compression in Language Models. In Proceedings of EMNLP 2023. (Oral)
- 8. Emily Cheng, Mathieu Rita, Thierry Poibeau. On the Correspondence between Compositionality and Imitation in Emergent Neural Communication. In Findings of ACL 2023.
- 9. Emily Cheng, Yen-Ling Kuo, Josefina Correa, Ignacio Cases, Boris Katz, and Andrei Barbu. Quantifying the Emergence of Symbolic Communication. In Proceedings of CogSci 2022.
- 10. Emily Cheng, Yen-Ling Kuo, Ignacio Cases, Boris Katz, and Andrei Barbu. Towards Modeling the Emergence of Symbolic Communication. In Proceedings of the ICRA-2021 Social Intelligence Workshop.

TALKS & PRESENTATIONS

September 2022-2027

January 2022 GPA: 4.7/5.0

Linearly Controlled Language Generation with Performative Guarantee Poster at MINT Workshop at NeurIPS	es December 2024 Vancouver	
Geometric Signatures of Compositionality Across a Language Model's Poster at NeurReps Workshop at NeurIPS	Lifetime December 2024 Vancouver	
Evidence from fMRI supports a two-phase abstraction process in langu Talk at RycoLab Seminar @ ETH Zurich Talk at Grammar and Cognition Group @ Universitat Pompeu Fabra Oral presentation at UniReps@NeurIPS	age models. November 2024 Zürich October 2024 Barcelona December 2024 Vancouver	
High-Dimensional Abstraction Phase in Language Transformers Talk at Area Science Park Talk at Infolab@MIT CSAIL group meeting Talk at Johns Hopkins Center for Speech and Language Processing	June 2024 Trieste, Italy June 2024 Cambridge, MA February 2025 Baltimore	
Bridging Information-Theoretic & Geometric Compression in LMs Oral presentation at Deep Learning Barcelona Symposium Talk at Zaslavsky lab (NYU Psychology) Oral presentation at EMNLP Talk at Swiss AI Lab (IDSIA) Talk at Evolution in Language (EviL) Seminar	December 2024 Barcelona August 2024 NYC December 2023 Singapore November 2023 Lugano September 2023 Online	
Interplay of functional & semantic aspects in shaping inflectional morp a case study on Romance languages Poster at Crosslinguistic Perspectives in Linguistics Conference (X-PPL) Talk at Evolution in Language (EviL) Seminar Poster at Mediterranean Morphology Meeting	hology: November 2023 Zürich December 2024 Online June 2025 Zadar, Croatia	
VDEDIENCE		

EXPERIENCE

Apple Machine Learning ResearchSummer 2025MLR InternCambridge, UK

· Calibrating prompt-based steering in generative models.

Universitat Pompeu Fabra

Doctoral thesis

• Explore relationship between linguistic structure, information-theoretic compression, and intrinsic dimensionality in neural language models.

Ecole Normale Supérieure, CNRS

Visiting researcher

- \cdot Supervised by Thierry Poibeau. Supported by Paris AI Research Institute.
- $\cdot\,$ Explore relationship between compositionality of emergent languages and ease of imitation learning.

MIT CSAIL Infolab

Master's Research

- \cdot Thesis: Understanding Symbolic Communication
- $\cdot\,$ Supervised by Boris Katz and Andrei Barbu.
- Characterized the transition from sub-symbolic to symbolic communication between human players, and later machine players via a communication game.

AWARDS & GRANTS

NeurIPS UniReps Workshop Best Abstract *Free registration to NeurIPS 2024*

UPF Department of Linguistics "Estada" Grant, € 3000 Competitive grant for research stay with MIT CSAIL Fall 2022-Barcelona

Spring-Fall 2022 Paris

Fall 2020 - Fall 2021 Cambridge, MA

> 2024 Vancouver, CA

 $\begin{array}{c} 2024\\Barcelona,\ ES\end{array}$

UPF Department of Linguistics Travel Grant , $\in 2000$ Competitive grant for travel to EMNLP conference	2023 Barcelona, ES
Brains, Minds, Machines Travel Grant Fully paid BMM Summer Course, $\approx 10\%$ acceptance	2023 Woods Hole, MA
Fulbright France Open Research Grant Semifinalist Meta-learning in low-resource multilingual generalization	2021 Paris, France
In collaboration with LATTICE at CNRS and École Normale Supérieure.	
TEACHING & SUPERVISION	
Brains Minds and Machines Summer Course Teaching Assistant	Summer 2024 Woods Hole, MA
Gave tutorials and supervised course projects.	
Madagascar ML Summer School Guest Lecture	Winter 2022
Introduction to Statistical Learning Theory.	
6.031 Software Construction Graduate Teaching Assistant	Fall 2020 Cambridge, MA
Held lab hours, graded assignments for students in intermediate Java software class.	
MIT-France Global Teaching Labs Instructor	January 2020 Grenoble, France
Taught middle, high, and preparatory school students for STEM outreach program.	
MIT Math Learning Center Teaching Assistant	Fall 2018 - Spring 2019 Cambridge, MA
Office hours for Differential Equations, Linear Algebra, Probability, Physics and Cal	culus.
MIT Math Department Grader	Fall 2017, Fall 2018 Cambridge, MA
Probability and Random Variables (18.600) and Statistics (18.650).	

ACTIVITIES & OUTREACH

Conference and workshop organization

CoNLL 2025 Publication Chair; COLT Symposium 2025;

Reviewer

NeurIPS, ICLR 2025; LanGame@NeurIPS workshop, UniReps@NeurIPS workshop, NeurIPS, ARR, ICLR 2024; EMNLP 2023; NeurIPS 2022; NeurIPS, ICLR 2021

INDUSTRY EXPERIENCE

Palantir Technologies

Software Engineering Intern

Developed insurance risk models using PySpark and built map visualization backend in Java, arcGIS

Two Sigma Investments

Summer 2019 New York, NY

Summer 2020

Remote

Quant Research Intern: News Team

Built alpha models in Python and Groovy to forecast equity and options returns with news data.

Virtu Financial <i>Quant Research Intern</i>		January 2019 New York, NY
Developed cross-asset market impact n	nodels in Python for cash equities execution	on.
Goldman Sachs Securities Quant Research Intern: Equ	ities Flow Vol, FICC SMM Execution Ser	Summer 2018 <i>New York, NY</i>
 Developed alpha models in Python to f Built an order fill model for trading sin 	forecast realized volatility for single stock nulation in Java	options
COURSEWORK		
Brains Minds and Machines Summer Course Project: Intrinsic Dimensionality of Brain Responses to Language		Summer 2023 Woods Hole, MA
Institute of Language, Communic	ation, and the Brain Summer Schoo	I Summer 2022 Marseille
Official Coursework 6.867 Machine Learning (G) 6.435 Bayesian Inference (G) 6.864 NLP (G) 6.884 Sensorimotor Learning (G)	 6.860 Statistical Learning Theory (G) 6.031 Software Construction 6.046 Design & Analysis of Algorithms 24.933 Semantics & Pragmatics (G) 	 6.337 Numerical Methods (G) 18.615 Stochastic Processes (G) 6.436 Probability Theory (G) 9.190 Comp. Linguistics (G)
SKILLS		

Computer Languages	Python, Java, C/C++, Julia
Software & Tools	Pandas/Numpy/Scipy, PyTorch, Git, Linux

LANGUAGES

English (native), Mandarin (fluent), French (C1), Spanish (C1)