

# Giorgio Giannone

[linkedin/giorgio-c-giannone](https://www.linkedin.com/in/giorgio-c-giannone)

ggiorgio@mit.edu — gigi@dtu.dk

github/georgosgeorgos

Principal Research Scientist with the AI Innovation Team at **Red Hat** and Research Affiliate at **MIT MechE**.

Specializing in Generative AI, Generative Design and Probabilistic Modeling, with a focus on **Inference-Time Scaling, Test-Time Adaptation, Vision-Language Alignment, and Few-Shot Generation**.

Leading research on Probabilistic Inference to develop efficient, grounded Foundation Models for data-constrained engineering domains.

## Experience

### Principal Research Scientist, Red Hat

- AI Innovation Team
  - **Research:** Probabilistic Inference for Vision and Language Models
  - **Research:** Inference-Time Scaling and Reasoning for LLMs (ICLR)
  - **Product:** its-hub Development and vLLM Gateway Integration
  - **Product:** Context Optimization for AgentOps

Boston, Massachusetts, USA

*June 2025 - Present*

### Research Affiliate, Massachusetts Institute of Technology

- DeCoDE Lab. Department of Mechanical Engineering
  - **Research:** Inference-Time Scaling for Constrained Generative Design
  - **Research:** Iterative Self-Training for CAD Program Synthesis

Cambridge, Massachusetts, USA

*Jan 2026 - Present*

### Applied Scientist, Amazon

- Home Innovation and GenAI Team
  - **Research:** Grounded Vision-Language Models
  - **Research:** Evaluation for Text-to-Image Models (CVPR)
  - **Product:** Detection and Ranking Algorithms for Amazon Visual Shopping
  - **Product:** Subject-Driven Generative Models for AI Creative Studio

Seattle, Washington, USA

*April 2024 - June 2025*

### Visiting Researcher, UCL Centre for Artificial Intelligence

- Host: David Barber
  - **Research:** Multi-Resolution Convolutional Models for Long Sequences (NeurIPS)
  - **Research:** Bayesian Inference for Language Models

London, UK

*Jan 2024 - March 2024*

### Researcher (PhD Intern), Microsoft Research

- ML and Statistics Group. Hosts: David Alvarez Melis, Nicolo Fusi
  - **Research:** Dynamic Vocabulary Augmentation for LLMs

Cambridge, Massachusetts, USA

*Jun 2023 - Sept 2023*

### Research Collaborator, MIT-IBM AI Lab

- Model Alignment Team. Host: Akash Srivastava
  - **Research:** Generative Models for Systems with Constraints (NeurIPS)
  - **Research:** Aligning Language Models with Negative Data
  - **Product:** Specialized Language Models for Enterprise Domains

Cambridge, Massachusetts, USA

*Jan 2023 - June 2023*

### Research Scientist (PhD Intern), IBM Research

- Accelerated Discovery Team. Hosts: Matteo Manica, Teodoro Laino
  - **Research:** Multitask Language Models for Text and Chemistry (ICML)
  - **Product:** Open-source library GT4SD for conditional generative models

Zurich, Switzerland

*Jun 2022 - Nov 2022*

### Applied Scientist (PhD Intern), Amazon Science

- Alexa Team. Hosts: Yunlong Jiao, Emine Yilmaz
  - **Research:** Domain Agnostic Subpopulation Generalisation

Cambridge & London, UK

*Jul 2021 - Oct 2021*

<b>Research Engineer, NNAISENSE</b>	Lugano, Switzerland <i>Jan 2019 - Jan 2020</i>
<ul style="list-style-type: none"> <li>◦ Deep Learning Team. Managers: Christian Osendorfer, Jonathan Masci           <ul style="list-style-type: none"> <li>– <b>Research:</b> Structured Latent Variable Models</li> <li>– <b>Product:</b> NeuralODE Algorithms for High-Range Event Camera Streams</li> </ul> </li> </ul>	
<b>Co-Founder, SecretAIry (formerly GAiA)</b>	Rome, Italy <i>July 2017 - Jan 2019</i>
<ul style="list-style-type: none"> <li>◦ Chatbots to enhance Workplace Communication           <ul style="list-style-type: none"> <li>– Selected among 100+ startups to join the EnLabs Incubator</li> </ul> </li> </ul>	

## Education

<b>PhD, Generative Machine Learning</b>	Technical University of Denmark, Lyngby, Denmark <i>June 2020 - Dec 2023</i>
<ul style="list-style-type: none"> <li>• Few-Shot Generative Models (ICML)</li> <li>• Multitask Language Models for Conditional Molecule Generation (ICML)</li> <li>• Diffusion Models for Generative Engineering Design and Topology Optimization (NeurIPS)</li> <li>• Thesis: Learning Generative Models with Limited Data           <ul style="list-style-type: none"> <li>– Supervisor: Ole Winther; Co-supervisor: Søren Hauberg</li> </ul> </li> </ul>	
<b>Visiting PhD Student, MIT School of Engineering</b>	Cambridge, Massachusetts, USA <i>Jan 2023 - Sept 2023</i>
<ul style="list-style-type: none"> <li>• Constrained Diffusion Models for Engineering Design (NeurIPS &amp; Patent)</li> <li>• Improving Generative Constraint Satisfaction using Invalid Designs (TMLR)</li> <li>• Evaluating Vision-Language Models for Engineering Tasks (Journal)</li> <li>• Research on LLM Agents for CAD design. Co-developer of <code>text2cad</code>.           <ul style="list-style-type: none"> <li>– Host: Faez Ahmed, DeCoDE Lab</li> </ul> </li> </ul>	
<b>Master's Degree, Data Science</b>	Sapienza University, Rome, Italy <i>Sept 2016 - Nov 2018</i>
<ul style="list-style-type: none"> <li>• Excellence Path &amp; Summa Cum Laude</li> <li>• Thesis: Multimodal Learning for Scene Understanding           <ul style="list-style-type: none"> <li>– Supervisor: Aris Anagnostopoulos; External Supervisor: Boris Chidlovskii</li> </ul> </li> </ul>	
<b>Visiting Graduate Student, NYU Tandon School of Engineering</b>	NYC, New York, USA <i>Sept 2017 - Jan 2018</i>
<ul style="list-style-type: none"> <li>• Visualization and Data Analytics Research Center. Host: Enrico Bertini           <ul style="list-style-type: none"> <li>– Built an interactive entity retrieval tool to investigate 10M documents</li> </ul> </li> </ul>	
<b>Master's Degree, Mechanical Engineering</b>	Sapienza University, Rome, Italy <i>Sept 2014 - Jan 2017</i>
<ul style="list-style-type: none"> <li>• Summa Cum Laude</li> <li>• Thesis: Bubble Dynamics in Turbulent Shear Flows           <ul style="list-style-type: none"> <li>– Supervisor: Carlo Massimo Casciola; Co-supervisor: Paolo Gualtieri</li> </ul> </li> </ul>	
<b>Bachelor's Degree, Mechanical Engineering</b>	Sapienza University, Rome, Italy <i>Sept 2009 - May 2014</i>
<ul style="list-style-type: none"> <li>• Thesis: Rapid Prototyping of Metallic Manufacturing</li> </ul>	

## Selected Publications & Patents

<b>Mitigating Premature Exploitation in Particle-based Monte Carlo for ITS</b>	under-review
<u>GIANNONE</u> , XU, NAYAK, AWHAD, SUDALAIRAJ, XU, SRIVASTAVA	<i>2025</i>
<b>Generative optimization models for machine learning</b>	US Patent (MIT & IBM)
<u>GIANNONE</u> , SRIVASTAVA, AHMED	<i>2025</i>
<b>Feedback-Driven Vision-Language Alignment</b>	under-review
<u>GIANNONE</u> , LI, FENG, PEREVODCHIKOV, CHEN, MARTINEZ	<i>2025</i>

<b>Be More Specific: Evaluating Object-centric Realism in Synthetic Images</b>	CVPR
LIANG, CORNEANU, FENG, <u>GIANNONE</u> , MARTINEZ	2025
<b>Evaluating Vision-Language Models for Engineering Design</b>	Springer Artificial Intelligence Review
PICARD, EDWARDS, DORIS, MANN, <u>GIANNONE</u> , ALAM, AHMED	2025
<b>Reparameterized Multi-Resolution Convolutions for Long Sequence Modelling</b>	NeurIPS
CUNNINGHAM, <u>GIANNONE</u> , ZHANG, DEISENROTH	2024
<b>Constraining Generative Models for Engineering Design with Negative Data</b>	TMLR
REGENWETTER, <u>GIANNONE</u> , SRIVASTAVA, GUTFREUND, AHMED	2024
<b>Aligning Optimization Trajectories with Diffusion Models</b>	NeurIPS
<u>GIANNONE</u> , SRIVASTAVA, WINTHER, AHMED	2023
<b>Unifying Molecular and Textual Representations via Multi-task LM</b>	ICML
CHRISTOFIDELLIS*, <u>GIANNONE</u> *, BORN, WINTHER, LAINO, MANICA	2023
<b>Accelerating Material Design with GT4SD</b>	Nature npj Computational Materials
<i>GT4SD Team (Core Contributor)</i>	2023
<b>Few-Shot Diffusion Models</b>	SBM@NeurIPS
<u>GIANNONE</u> , NIELSEN, WINTHER	2022
<b>SCHA-VAE: Hierarchical Context Aggregation for Few-Shot Generation</b>	ICML
<u>GIANNONE</u> , WINTHER	2022
<b>Method and apparatus for semantic segmentation and depth completion</b>	US Patent (NAVER)
CHIDLOVSKII, <u>GIANNONE</u>	2022

## Projects & Open Source

<b>its-hub: A Python library for inference-time scaling</b>	2025
– Contributor.	
– Inference-Time Scaling for Language Models.	
– Focus on Mathematical Reasoning.	
– Contributed Entropic Particle Filtering algorithms and new benchmark.	
<b>Text2CAD: Democratizing Engineering Design. Prompt by Prompt.</b>	2023
– Co-Lead.	
– DesignX. Team of engineers and researchers based at MIT and Caltech.	
– Generative tool that allows users to create CAD models using natural language prompts.	
– The tool is designed to be user-friendly and accessible to non-experts, enabling a wide range of users to quickly create complex CAD models without the need for specialized training.	
<b>GT4SD: Generative Toolkit for Scientific Discovery</b>	2022
– Core Contributor.	
– Library leveraging conditional generative models for accelerated discovery.	
– Work on Diffusion Models for images and 3D molecule conformation. The GFlowNet framework. Property Prediction module. Public Hub for model upload. Training Pipelines. Documentation. Tutorials. Testing. CI/CD. Server and Client API. Docker Images for CPU and GPU.	

## Grants & Awards

<b>GPU Grant, LUMI-G, EuroHPC</b>	Copenhagen, Denmark
PI, Efficient Pre-training of Large Generative Models for Constrained Design	Nov 2023
<b>Grant, Otto Møensted's Foundation</b>	Copenhagen, Denmark
Research Grant	Dec 2022
<b>Grant, Independent Research Fund Denmark</b>	Lyngby, Denmark
	Jun 2020

DFF PhD Grant

**Grant, Perception as Generative Reasoning**

Awarded complimentary NeurIPS registration by DeepMind

NeurIPS 2019

*Oct 2019*

**Grant, Pi School**

Full Tuition for the School of AI (3% acceptance rate)

Rome, Italy

*Oct 2018*

**Certificate of Award, Tsinghua University**

Prize for Outstanding Accomplishments, Deep Learning Summer School

Beijing, China

*Aug 2018*

**1st Pick, Excellence Path, Master's Degree, Data Science**

Admission based on the First year's Academic Performance

Rome, Italy

*Mar 2018*

Participation in the School for Advanced Studies

## Academic Service

### Reviewer

Conference: ICML19, ICCV19, AAAI20, ICML21 (top 10%), AISTATS21, ICML22, NeurIPS22, CVPR23, NeurIPS23, ICML24, ICLR25, CVPR25, NeurIPS25, ICLR26

Journal: TPAMI, TMLR

Workshop: NeurIPS-IBW20, NeurIPS-MetaLearn21, ICML-DeployableGenAI23, ACL-LanguageMolecules24

### Teaching

Teaching: Deep Learning (DTU 02456), Bayesian Machine Learning (DTU 02477), Advanced Machine Learning (DTU 02460)

Supervision: two special courses (9 months), two master's thesis (6+6 months), 18 final projects

### Volunteering

PAISS18, NeurIPS18, ELLIS Unit Copenhagen, MLLS

## Skills

### Languages

- Python (proficient); R, Matlab (good knowledge); C, Java, JavaScript (basic knowledge)

### Research

- Accelerate, HF Transformers, LaTeX, NLTK, OpenCV, PyTorch, SpaCy, TensorFlow, verl

### Software

- AWS, CVX, Docker/podman, FastAPI, Git, GitHub Actions, Gradio, Linux, MinIO, MongoDB, MySQL, Travis, vLLM, LangGraph, LangFlow, Langfuse, Cline, Cursor, OpenRouter