

# VEDANT PURI

## OBJECTIVE

vedantpuri@cmu.edu | <https://vpuri3.github.io/>

PhD candidate graduating August 2026 seeking a technical role in machine learning research. Experience developing scalable-transformers and custom GPU kernels for large-scale training. Background in computational fluid dynamics.

## EDUCATION

<b>Carnegie Mellon University</b> <i>Ph.D Mechanical Engineering.</i> Advisors: Levent Burak Kara, Yongjie Jessica Zhang Proposed thesis: <i>Neural representations for computational physics: from reduced order modeling to transformers</i>	Jan 2022–Present
<b>University of Illinois Urbana-Champaign</b> <i>B.S. Engineering Mechanics, B.S. Mathematics.</i>	2015–2019

## SELECTED RESEARCH CONTRIBUTIONS

<b>FLARE: Fast Low-Rank Attention Routing Engine</b>   <i>Efficient attention architectures</i>	2025
• Derived a flexible low-rank reformulation of self-attention via latent routing.	
• Reduced quadratic complexity of global communication in self-attention to linear complexity.	
• Demonstrated scaling to 1M tokens on a single H100 GPU, attaining over 200× speedup over vanilla self-attention.	
<b>FLARE for Language Modeling (Ongoing dissertation work)</b>   <i>Decoder-only architectures</i>	2025–Present
• Extending FLARE to decoder-only next-token prediction with causal attention.	
• Enabling adaptive attention state size to control memory and compute during training and inference.	
• Implementing fused Triton kernels for causal training, prefill, and decode.	
<b>Equation-based PDE modeling with neural manifolds</b>   <i>Hybrid data + physics methods</i>	2024
• Introduced smooth neural fields as nonlinear spatial ansatz functions in equation-based reduced-order modeling.	
• Retained physics-based Galerkin time evolution while learning expressive low-dimensional representations.	
• Attained 199× speedup over full-order simulations in transport-dominated regimes.	

## EXPERIENCE

<b>Carnegie Mellon University</b>   <i>Research Assistant</i>	Jan 2022–Present
• Wrote solvers for phase field simulations of lithium dendritic growth in solid-state batteries.	
• Wrote fluids dynamics solvers for turbulence closure modeling with differentiable physics.	
<b>Julia Computing</b>   <i>Intern Engineer</i>	Apr 2021–Nov 2021
• Built numerical solvers for scientific machine learning ecosystem in Julia.	
<b>Carnegie Mellon University</b>   <i>Research Assistant</i>	Sep 2020–Jan 2021
• Developed differentiable geometry representations and meshing algorithms.	
<b>Argonne National Laboratory</b>   <i>Research Assistant</i>	Mar 2020–Sep 2020
• Executed large-scale simulations of turbulent airflow over urban landscapes on supercomputers.	
<b>Argonne National Laboratory</b>   <i>Research Assistant</i>	May 2018–Jul 2018
• Executed high-fidelity fluid dynamics simulations and analyzed turbulent statistics for closure modeling.	
<b>National Center for Supercomputing Applications</b>   <i>Intern</i>	Sep 2017–May 2018
• Numerical simulation of spacetime metric for gravitational wave simulations.	

## PUBLICATIONS

- Puri, V. et al., *FLARE Decoder: Low-rank attention routing for causal language modeling* (In preparation).  
Puri, V. et al., *FLARE: Fast Low-Rank Attention Routing Engine*. arXiv:2508.12594 (2025) (Under review).  
Puri, V. et al., *Reduced order modeling with smooth neural fields*. JCP 2025, doi:10.1016/j.jcp.2025.113957.  
S, V., Puri, V., et al., *Closure modeling for Burgers' turbulence*. MLST 2023, doi:10.1088/2632-2153/acb19c.

## AWARDS

<b>World Conference on Computational Mechanics</b>   <i>Best poster in fluid dynamics</i>	2024
<b>University of Illinois</b>   <i>Theoretical and Applied Mechanics Merit Award</i>	2019

## TECHNICAL SKILLS

- ML Architectures:** Efficient attention, causal transformers, memory-efficient prefill/decode, neural operators.  
**ML Systems:** PyTorch, distributed/mixed-precision training, distributed training, causal language modeling.  
**GPU & Systems:** Triton kernel development, FlashAttention-style block algorithms, online softmax, Nsight profiling.  
**Scientific Foundations:** Numerical linear algebra, PDE-constrained modeling, finite elements, spectral methods.  
**Programming Languages:** Python, Julia, C, Fortran 77, MATLAB, UNIX, L<sup>A</sup>T<sub>E</sub>X.