## Prasanna Biswas

AI Software Solutions Engineer at Intel Corporation

Contact

## Work experience

present	AI Software Solutions Engineer	Email prasanna.biswas140gmail.com
↑ Jan 2024	Kernels, Falcon Shores, Intel Corporation	Ph en e
	• Developing <b>high-performance kernels</b> with dynamic shape support for Intel's next-gen GPU using SYCL, optimizing la-	(+91) 9922365239
	<ul> <li>tency, memory bandwidth, I/O access &amp; compute utilization.</li> <li>Programmed an efficient cumsum kernel, achieving 2x perf</li> </ul>	in Profile /in/prasanna-biswas
	<ul> <li>improvement over IPEX eager mode implementation.</li> <li>Designed and implemented complex operations like TopK and media operators such as Brightness and Contrast as graphs in C++ using MLIB types and attributes enabling</li> </ul>	Portfolio prasannabiswas-iitb.github.io
	efficient GPU execution.	M Tech, Thesis
	<ul> <li>Innovated a novel machine learning algorithm combining VAEs and Diffusion Models for NLP and CV.</li> <li>Co-authored two papers; one submitted to CVR 2025 and second submitted to IEEE CONNECT-2025.</li> </ul>	<b>Computational Model to Understand</b> <b>Emotions in Sarcasm</b> Created the 'emo-UStARD' dataset by annotating 'MUStARD' with 8 emo-
Jan 2024 🌘	Senior ML Engineer	• tions, arousal, and valence.
↑ Dec 2022	ML Applications, Cloud Al100, Qualcomm CR&D	Conducted experiments, observing an
	• Spearheaded ONNX optimizations on Qualcomm's AI100 accelerator, achieving an 8.5% performance boost for large	<b>18% increase</b> in accuracy across various aspects of textual modality.
	<ul> <li>language models (LLMs) like ChatGLM2-6B through node- fusions, graph simplifications.</li> <li>Enhanced GPT model efficiency by 2x through caching Key-</li> </ul>	Technical Blogging & Content Creation
	<ul><li>Value matrices and minimizing DDR reads/writes.</li><li>Designed a Graph Neural Network algorithm to enhance</li></ul>	Fechnical Blogs     GPUs and CUDA Programming
	<ul> <li>compiler efficiency, resulting in a filed patent.</li> <li>Led a three-member team in optimizing and deploying the top 120 models from Hugging Face library.</li> </ul>	YouTube Channel Co-Owner & W Python Instructor
Nov 2022	ML Engineer	Successfully manage a channel with 1.5k+ subscribers.
↑ Nov 2020	ML Applications, Cloud Al100, Qualcomm CR&D	Technologies
	<ul> <li>Engineered software modules in C++ &amp; Python.</li> <li>Introduced auto-detection of post-processing in CV models, replacing them with ABP &amp; NMS optimized kernels for 80% improvement in quantization accuracy.</li> <li>Achieved a 28.2% perf improvement for (BERT and variants) through Graphcore's packing strategy.</li> <li>Enhanced operator support in the GLOW compiler for the Cloud AI100 SDK.</li> </ul>	<ul> <li>Programming:</li> <li>Python, C++</li> <li>GPU: SYCL(DPC++), CUDA</li> <li>Machine Learning Frameworks:</li> <li>PyTorch</li> <li>ONNX, ONNX Runtime</li> <li>ML Domain &amp; Techniques:</li> </ul>
	Patent and Publications	<ul> <li>NLP, CV</li> <li>Graph Optimization, GNN</li> <li>Quantization, Pruning, Node Fusion</li> </ul>
Mar 2025	Efficient Deep Learning Model Architecture for Emergence of Machine Style Calligraphy	• GPU Optimization
	IEEE CONNECT-2025 Conference (Submitted)	• Git, Docker
Dec 2024	Generating Machine-Style Handwriting: A Diffusion-based latent generation with VAE decoding	• GLOW (Machine Learning Compiler)
	CVR 2025 Conference (Accepted for presentation)	Education
Jun 2023	Pre-Processing For Deep Neural Network Compilation Us- ing Graph Neural Networks	M Tech, 2020 • IIT Bombay
	USPTO: 18/330,253 and 18/500,014 (Pending)	CP1: 8.43/10
Jun 2018	Home Automation Using Panoramic Image Using IoT	<b>b</b> Tech, 2018 • VESIT, Mumbai CPL 0.07/10
1	ruoiisnea in: 2018 ICKIEECE	OF1: 9.07/10