



Hao Fang¹, Rajeev B. Botadra¹, ChiJui Chen⁵, Leo Scholl¹, Trung Le¹, Ryan Canfield², Eli Shlizerman³, Bo-Cheng Lai⁵, Amy Orsborn¹, Scott Hauck¹ ¹Electrical & Computer Engineering, University of Washington ²Bioengineering, University of Washington ⁵Electrical Engineering, National Yang Ming Chiao Tung University

- milliseconds
- impose large latencies



FPGA Accelerated Brain-Computer Interfaces A Co-design for Future Integrated Systems







Results

FPGA acceleration enables submillisecond inference with LFADS, a speedup of 1650x, 118x versus a CPU, GPU

Architecture Comparison		
itecture	Inference Latency (ms)	Power (W)
n 16-Core CPU	257.516 ± 56.693	84.345
FX3090 GPU	18.551 ± 2.995	87.393
J55C FPGA	0.156 ± 0.002	27.867

Hardware resources are not fully utilized, can support additional compute still (i.e larger models, additional postprocessing)

Resource Usage on U55C (8-bit LFADS)		
esource	Utilization	
LUTs	13.90%	
ip Flops	6.15%	
BRAM	6.03%	
DSPs	18.44%	

Future Work

Our co-designed system serves as a silicon testbed for future real-time experiments Extending compatibility to various recording modalities (e.g., EEG) b) Exploring advanced neural decoder models (e.g., transformer architectures) Investigating neuroplasticity^[2]

References

[1] J. Robinson et al., "An application-based taxonomy for braincomputer interfaces," Nat. Biomed. Eng, 2024 [2] A. L. Orsborn et al., "Closed-loop decoder adaptation shapes neural plasticity for skillful neuroprosthetic control," Neuron, 2014. [3] C. Pandarinath et al., "Inferring single-trial neural population dynamics using sequential auto-encoders," Nat. Methods, vol. 15, no. 10, pp.805–815, 2018.