# Biological Blueprints for Human Inspired AI

Stanford, April 7 & 9, 2020

# Biological Blueprints for Human Inspired AI‡

What does it mean for a brain to perform computations?†

# There are Many Approaches to Studying the Brain



*Blind monks examining an elephant*, an Ukiyo-e print by Hanabusa Itchō (1652–1724).
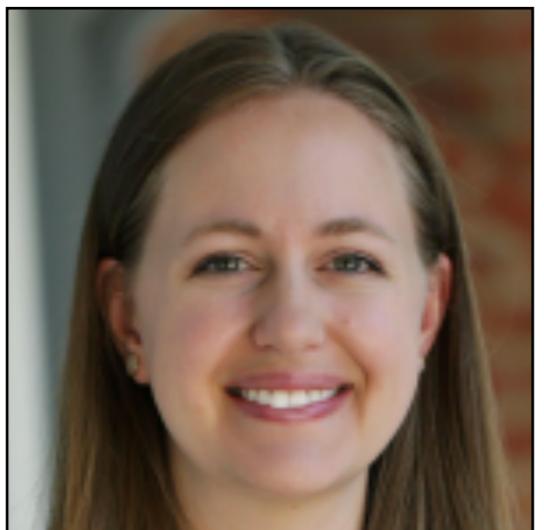
Adam Marblestone

Brenden Lake

Jessica Hamrick

Jill Leutgeb

Lisa Giocomo

Loren Frank

Matt Botvinick
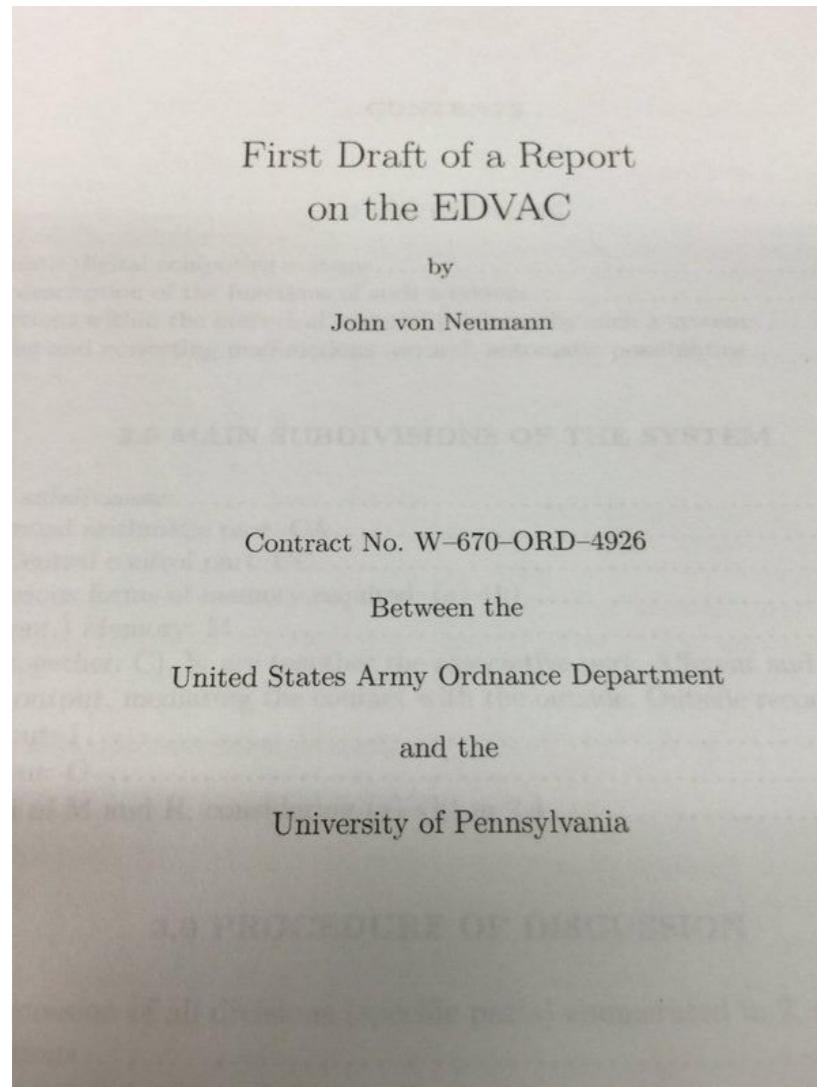
Michael Frank

Oriol Vinyals
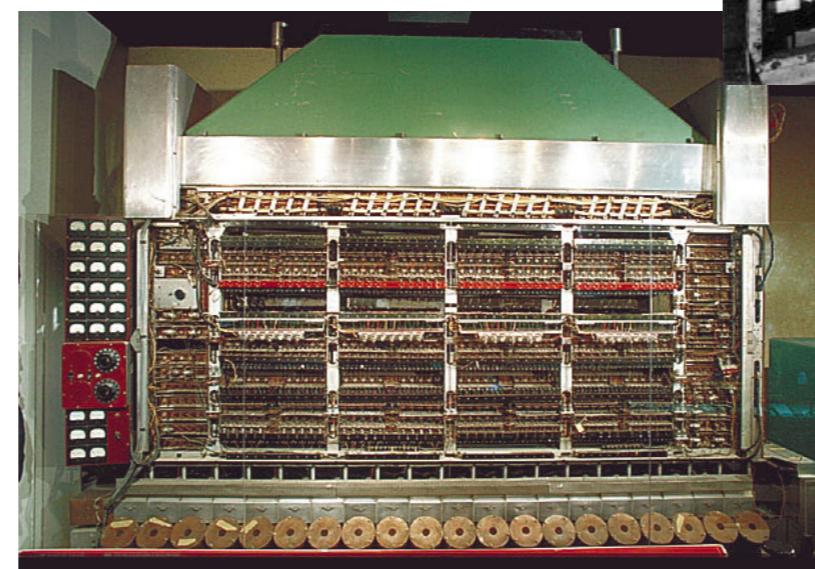
Randy O'Reilly

Peter Battaglia
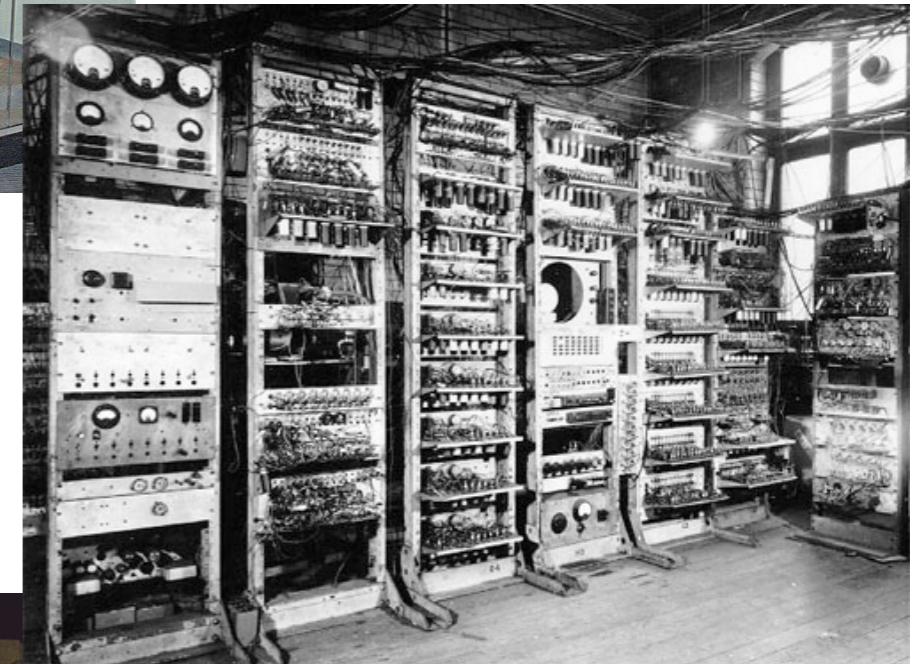
Vivek Jayaraman

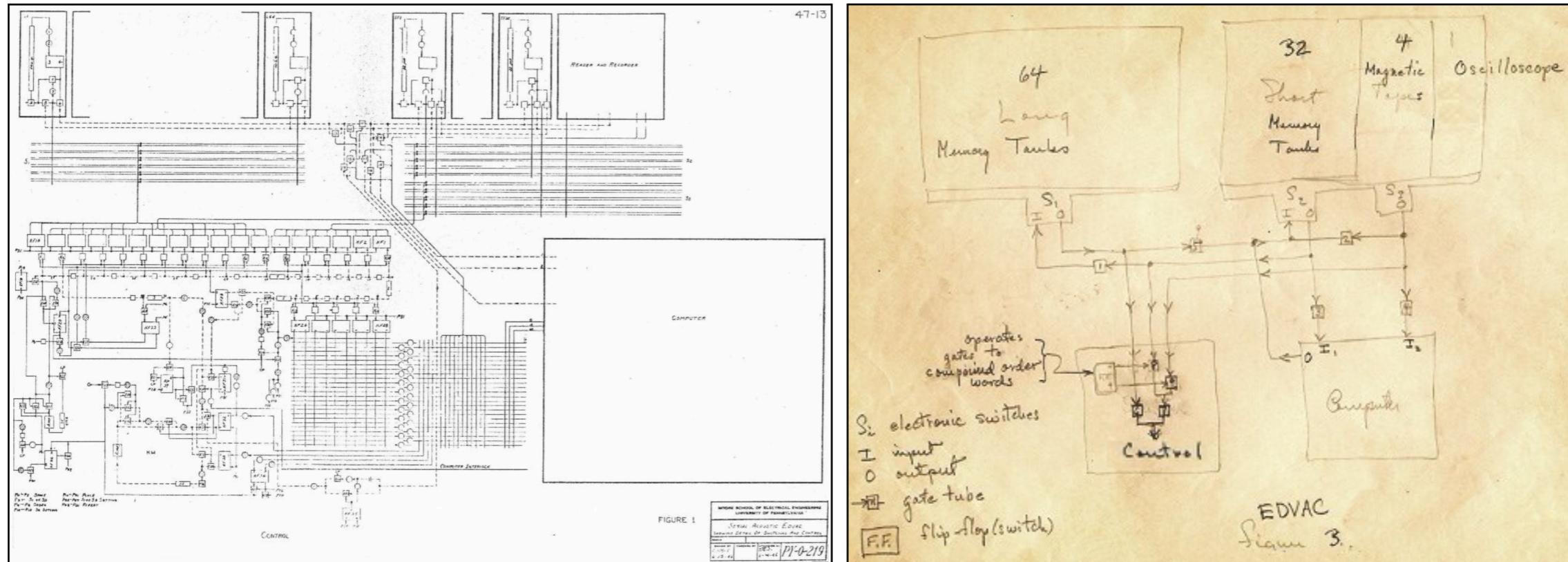# Birth of the Modern Computer:  The von Neumann Architecture



ENIAC



EDVAC

First Draft of a Report
on the EDVAC

by

John von Neumann

Contract No. W–670–ORD–4926

Between the

United States Army Ordnance Department

and the

University of Pennsylvania

IAS

Alice Wang, Benton Calhoun, and Anantha Chandrakasan. *Sub-Threshold Design for Ultra Low-Power Systems.* Springer-Verlag, 2006.

What would a human-brain analog of von Neumann's Architecture look like?

# Birth of the Modern Computer:  The von Neumann Architecture





**Von Neumann Basic Structure**





Josh Merel, Matthew Botvinick, and Gregory Wayne. Hierarchical motor control in mammals and machines. Nature Communications, 10(1):5489, December, 2019.

# Shannon, Turing, Gödel & von Neumann: <u>The Digital Abstraction</u>†



Alan Turing
Kurt Gödel

Claude Shannon

Hennessy & Patterson

John von Neumann

# Connections:
# It's the <u>Network Dummy</u>†

Is there an analog of the Digital Abstraction for modeling the human brain?

# Artificial Intelligence and the History of Connectionism

SYMBOLIC

CONNECTIONIST

Fodor & Pylyshyn (LoT)

Rumelhart, McClelland & Hinton (PDP)

COMBINATORIAL
COMPOSITIONAL
SERIAL

DISTRIBUTED
DIFFERENTIABLE
PARALLEL

Jerry A. Fodor. *The Language of Thought*. Harvard University Press, Cambridge, MA, 1975.

Jerry A. Fodor and Zenon W. Pylyshyn. Connectionism and cognitive architecture. *Cognition*, 28(1-2):3-71, 1988.

G. E. Hinton, J. L. McClelland, and D. E. Rumelhart. Chapter 3: Distributed Representations. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing, Explorations in the Microstructure of Cognition: Foundations*. MIT Press, Cambridge, MA, 1986.

Randall C. O'Reilly, Alex A. Petrov, Jonathan D. Cohen, Christian J. Lebiere, Seth A. Herd, and Trent Kriete. How limited systematicity emerges: A computational cognitive neuroscience approach. In Paco Calvo and John Symons, editors, *The Architecture of Cognition*, pages 191-224. MIT Press, Cambridge, Massachusetts, 2014.

Ortwin Bock. Santiago Ramón y Cajal, Camillo Golgi, Edward Schäfer and the Neuron Doctrine. *Endeavour*. 37(4):228–234, 2013.

Garcia-Lopez, Garcia-Marin, Miguel Freire. The histological slides and drawings of Cajal. *Frontiers in Neuroanatomy*, 4:1-16, 2010.

# Zeiss Multi-Beam Scanning Electron Microscope



High throughput multi-beam EM imaging with up to 91 parallel beams operating at 2 terapixels per hour with 3.5 nm resolution or better

Thomas Dean, Biafra Ahanonu, Mainak Chowdhury, Anjali Datta, Andre Esteva, Daniel Eth, Nobie Redmon, Oleg Rumyantsev, and Ysis Tarter.  On the technology prospects and investment opportunities for scalable neuroscience.  *CoRR*,  arXiv:1307.7302,  2013.

# Zebra Finch

# Drosophila Melangaster

C. Shan Xu, Michal Januszewski, Zhiyuan Lu, Shin-ya Takemura, Kenneth J. Hayworth, Patricia K. Rivlin, Vivek Jayaraman, […], Gerald M. Rubin, Harald F. Hess, Louis K. Scheffer, Viren Jain, and Stephen M. Plaza. A connectome of the adult drosophila central brain. *bioRxiv*, 2020. VIDEO

# Drosophila Melangaster



Sophie Aimon,  Takeo Katsuki,  Logan Grosenick,  Michael Broxton,  Karl Deisseroth,  and  Ralph J. Greenspan. Activity sources from fast large-scale brain recordings in adult drosophila. *bioRxiv*, 2015.

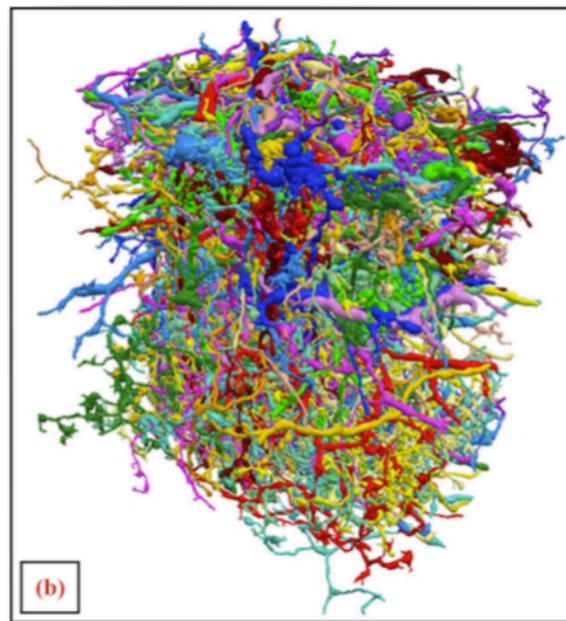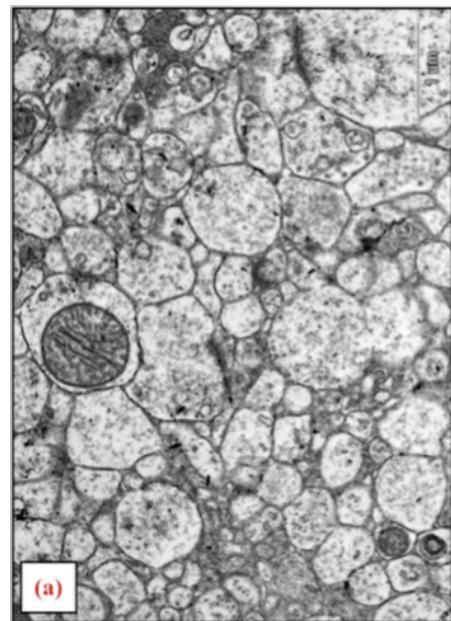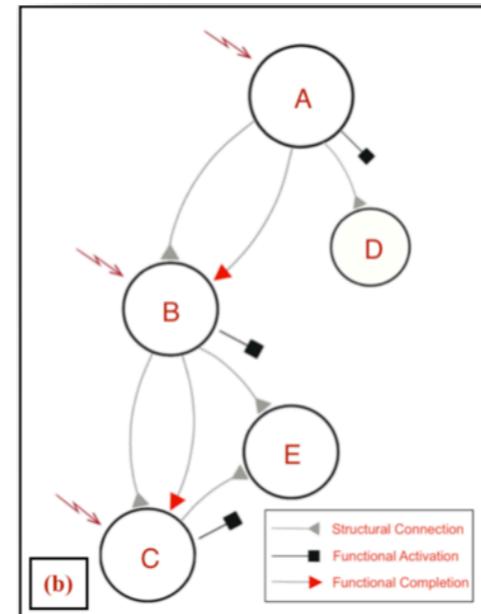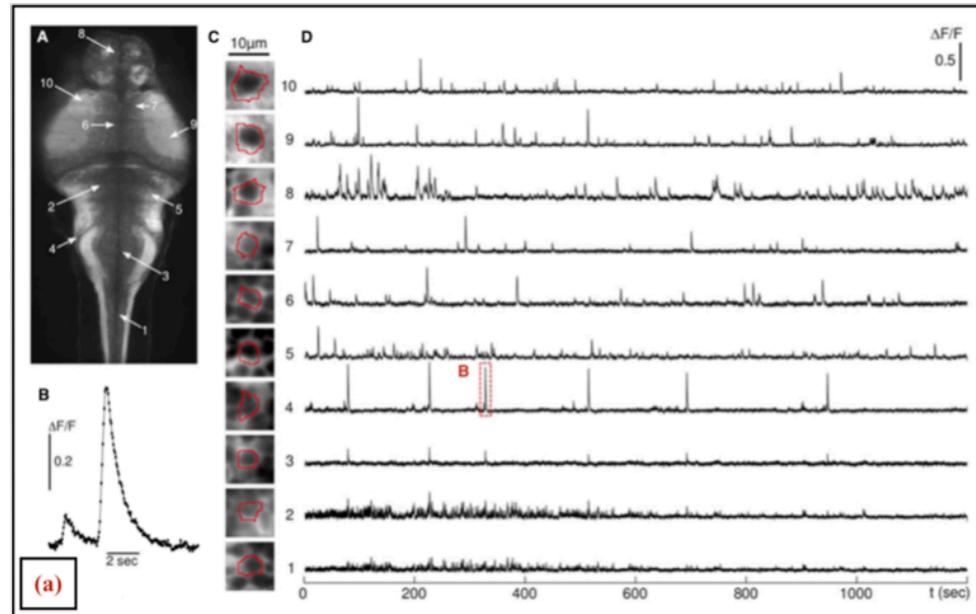# Mesoscale Modeling

1.(a) serial section sample
1.(b) dense reconstruction
1.(c) synaptic elaborations

2.(a) dense 2PE recording
2.(b) synaptic transmission

3.(a) 3D embedding space
3.(b) KD tree NN database
3.(c) configurable network
3.(d) basis function filters
3.(e) local and global loss

Thomas Dean. Inferring mesoscale models of neural computation. *CoRR*, arXiv:1710.05183, 2017.