# RECORD SIMULATION OF THE FULL-DENSITY SPIKING POTJANS-DIESMANN-MICROCIRCUIT MODEL ON THE IBM NEURAL SUPERCOMPUTER INC 3000

Arne Heittmann[1], Georgia Psychou[1], Guido Trensch[2], Charles E. Cox[3], Winfried W. Wilcke[3], Markus Diesmann[4] and Tobias G. Noll[1]
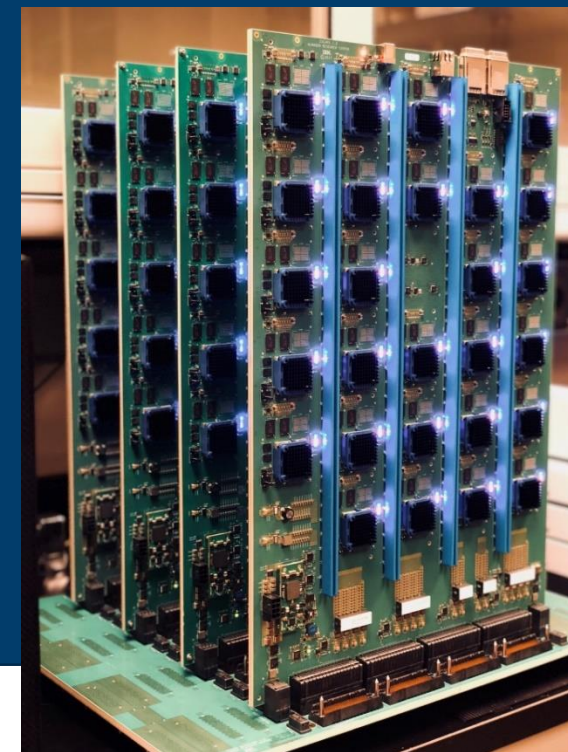
March 30 2022

ACA- towards multi-scale natural-density Neuromorphic Computing

[1]JARA-Institute Green IT (PGI-10), Jülich Research Centre, D-52425 Jülich, Germany
[2]Simulation & Data Lab Neuroscience, Jülich Supercomputing Centre Jülich Research Centre, D-52425 Jülich, Germany
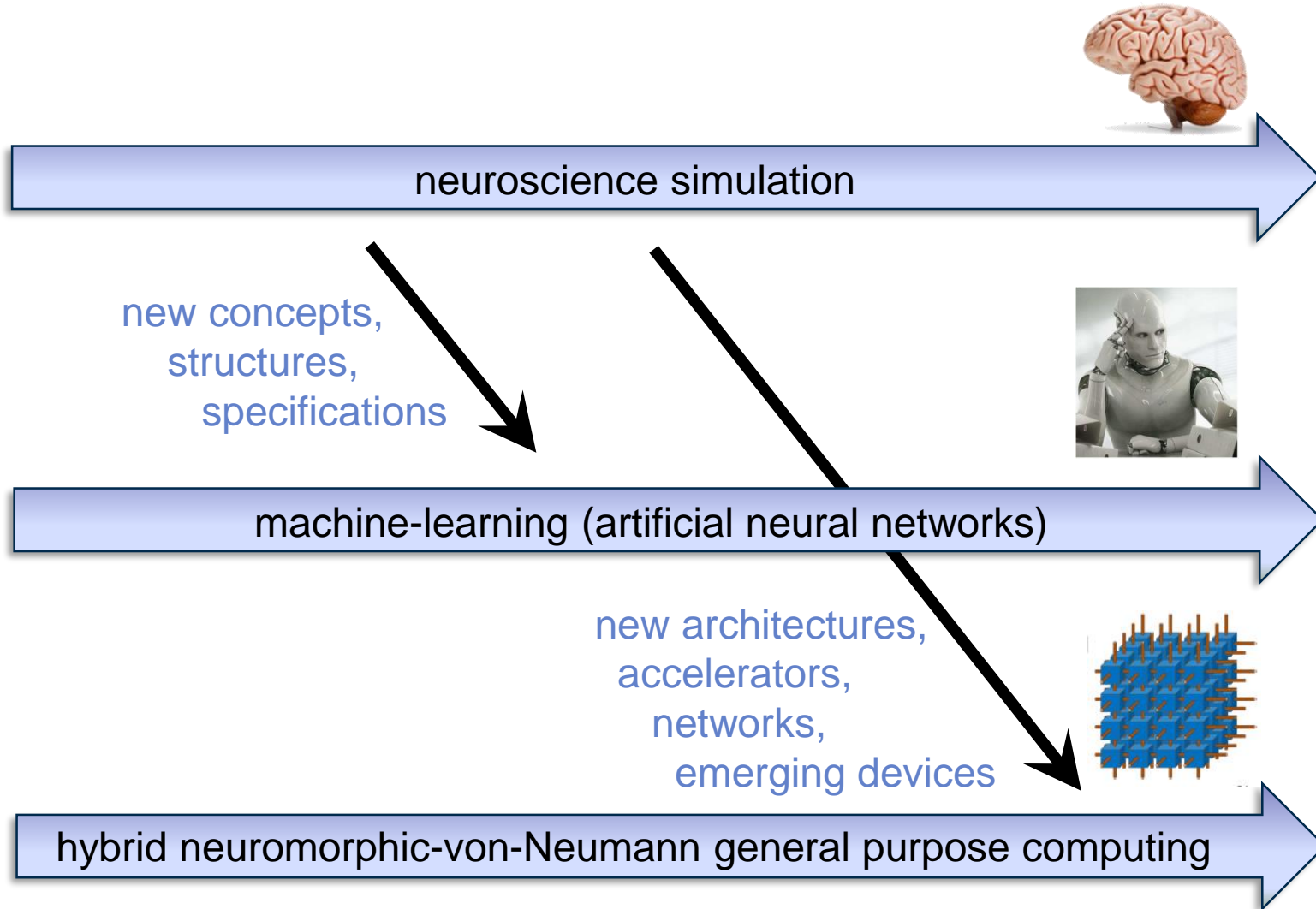[3]IBM Research Division, Almaden Research Center, San Jose, CA 9512
[4]Institute of Neuroscience and Medicine (INM-6)

# NEUROMORPHIC COMPUTING

- towards advanced general purpose computing architectures

**Neuroscience**

- dynamics of natural brains

- function of natural brains

neuroscience simulation

- learning, plasticity, development

new concepts, structures, specifications

**Artificial General Intelligence (AGI)**

- deep learning using few samples

machine-learning (artificial neural networks)

- ability for contextual adaptation

new architectures, accelerators, networks, emerging devices

- ability to explain descisions in natural language

hybrid neuromorphic-von-Neumann general purpose computing

# ACA – ADVANCED COMPUTING ARCHITECTURES

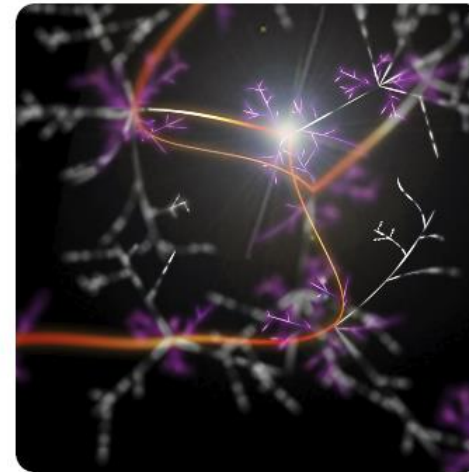**towards multi-scale natural-density neuromorphic computing**

- Pilot project preparing a long-term research initiative in the application area of *Neuroscience Simulation*

- Specification of a neuromorphic computing architecture for accelerating simulation experiments



Requirements, Validation & Benchmarking

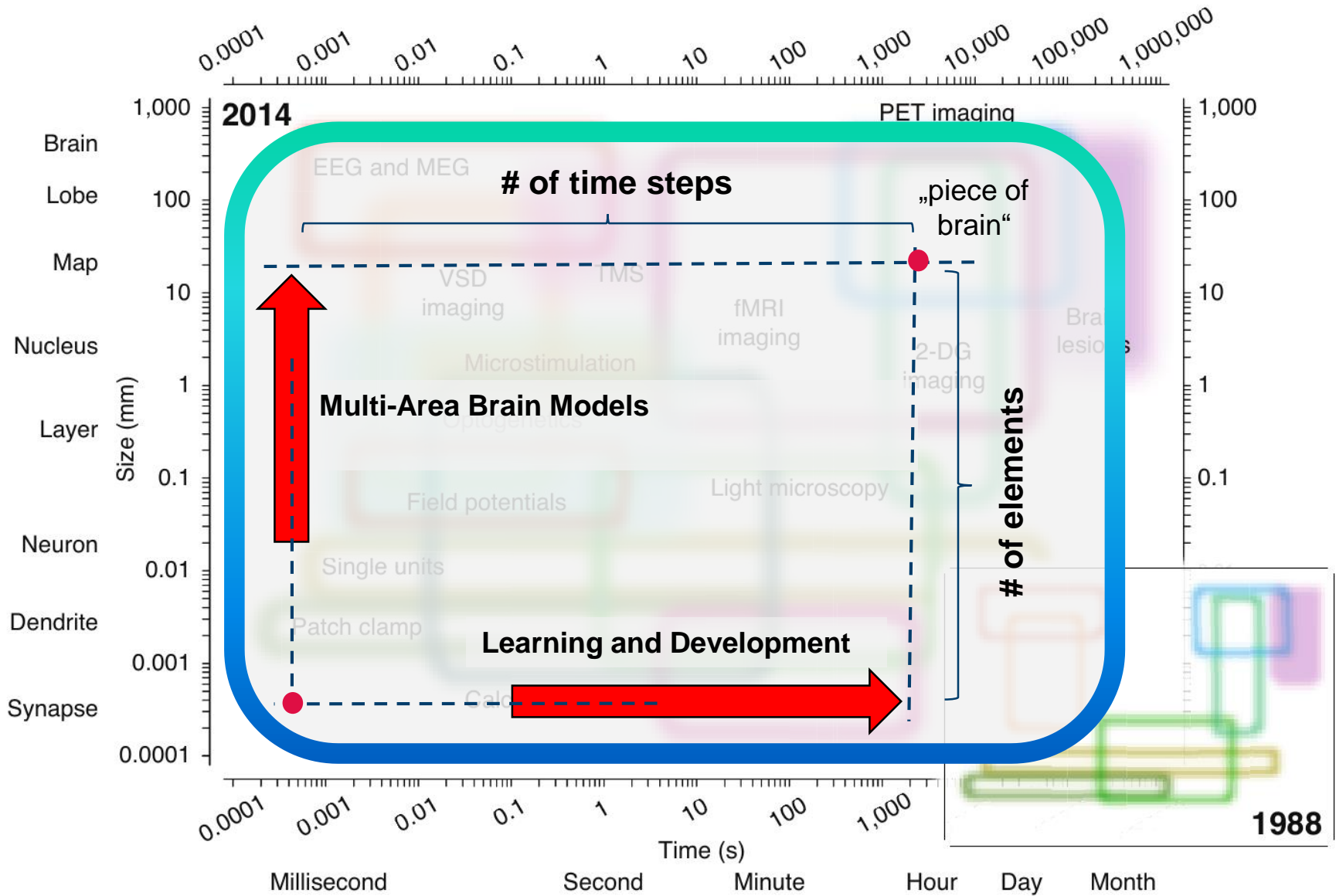

System Definition, Integration & Operation



Network Connectivity & Communication



Accelerated Numerics

# SCALES OF BRAIN ORGANIZATION



Sejnowski et al. (2014) Nature Neuroscience

# Trends in Neuroscience Simulation Experiments

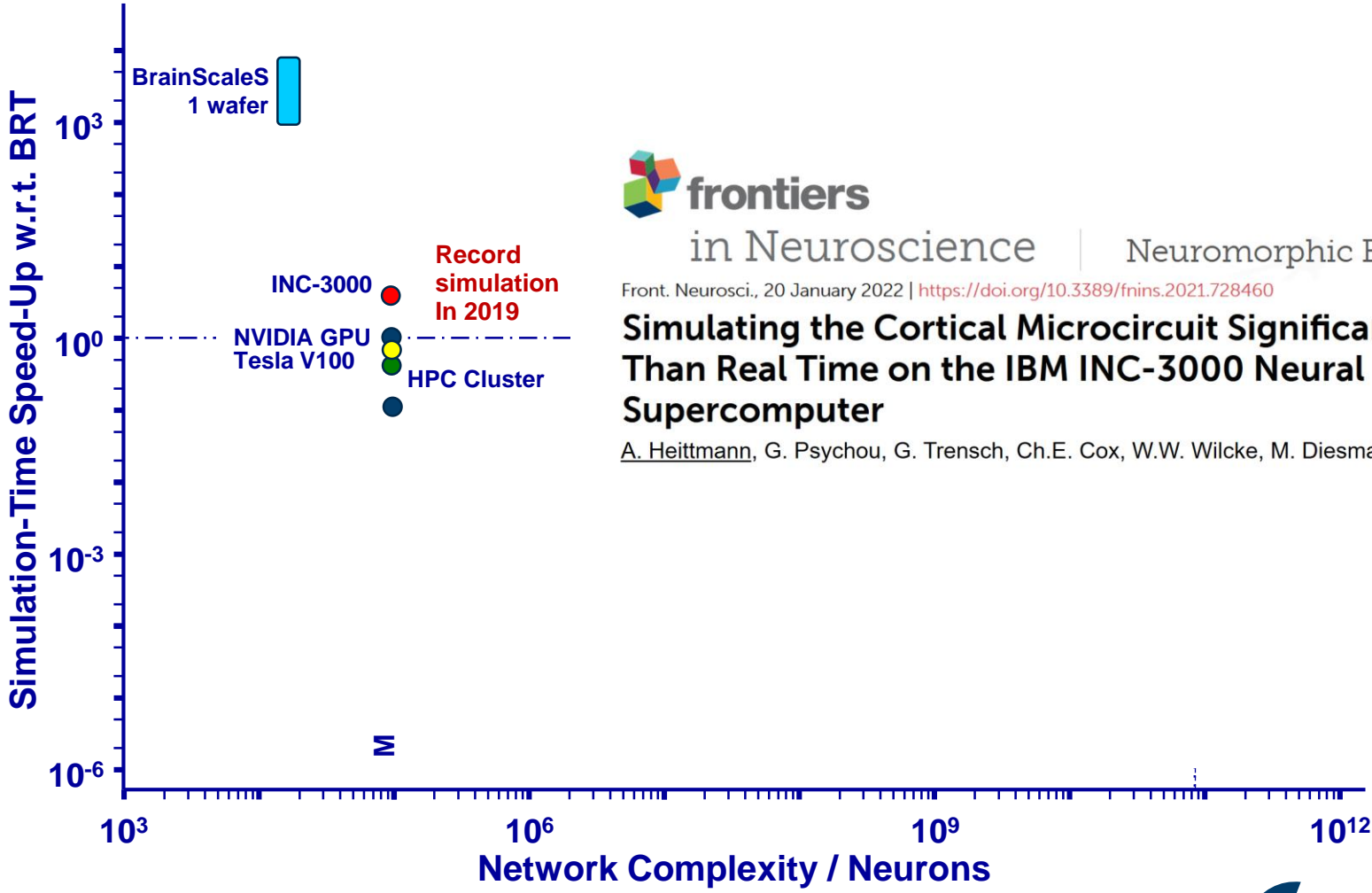BrainScaleS,
Univ. Heidelberg

IBM INC 3000
FZ-Jülich

JUQUEEN
FZ-Jülich

K-Computer Riken

**State-of-the-Art Simulation Time** (0.1-ms time grid)



**BrainScaleS 1 wafer**

**INC-3000**
**Record simulation In 2019**

**NVIDIA GPU Tesla V100**
**HPC Cluster**

**M**

Simulation-Time Speed-Up w.r.t. BRT

$10^3$
$10^0$
$10^{-3}$
$10^{-6}$

$10^3$ $10^6$ $10^9$ $10^{12}$

**Network Complexity / Neurons**

frontiers in Neuroscience | Neuromorphic Engineering

Front. Neurosci., 20 January 2022 | https://doi.org/10.3389/fnins.2021.728460

Simulating the Cortical Microcircuit Significantly Faster Than Real Time on the IBM INC-3000 Neural Supercomputer

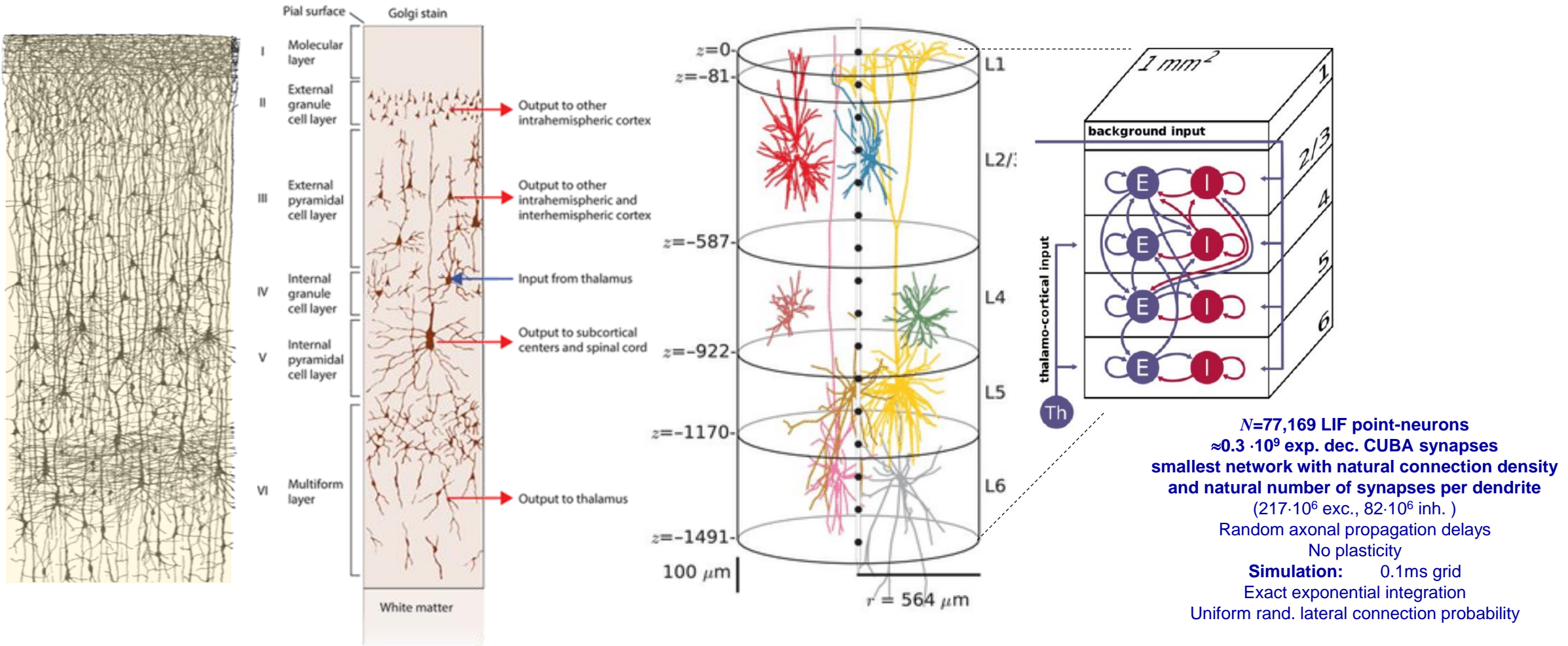A. Heittmann, G. Psychou, G. Trensch, Ch.E. Cox, W.W. Wilcke, M. Diesmann, and T.G. Noll

JÜLICH
Forschungszentrum

# THE CORTICAL MICROCIRCUIT

## The Cell-Type Specific Cortical Microcircuit: Relating Structure and Activity in a Full-Scale Spiking Network Model

Tobias C. Potjans[1,2,3] and Markus Diesmann[1,2,4,5]

Source: Tony Mosconi, Victoria Graham:
*Neuroscience for Rehabilitation*
Copyright © McGraw-Hill Education. All ri[...]

$N$=77,169 LIF point-neurons
≈0.3 ·10⁹ exp. dec. CUBA synapses
smallest network with natural connection density
and natural number of synapses per dendrite
(217·10⁶ exc., 82·10⁶ inh. )
Random axonal propagation delays
No plasticity
**Simulation:** 0.1ms grid
Exact exponential integration
Uniform rand. lateral connection probability

# THE IBM INC 3000 NEURAL SUPERCOMPUTER

- Originally developed for the IBM AGI (Artificial General Intelligence) Project (IBM Research, Almaden, CA)

- Platform for development and evaluation of prototypical circuit and architecture concepts for ACA

**INC 3000**



- INC board
- 27 x Xilinx XC 7Z045 SOC

- XC 7Z045 SOC
- 200 k LUTs
- 400 k FFs
- 900 DSPs
- 19.2 Mb BRAM
- 12.5 Gb/s GTX tranceivers
- 2 x ARM A9

- FPGA logic for Simulator, Router, and Configuration

- 16 x INC board
- 432 SoC Nodes

**Communication Network**



**Communication Network Topology:**

12 x 12 x 3   3D-node mesh

**Cross-sectional bandwidth:**

B = 450 Gb/s

**Worst case packet-path (betw. A and B)**

Latency T=24 μs

# XILINX 7Z045 SOC: THE COMPUTE NODE



**1 GB SRRAM**
external

processing system (PS)

ARM A9 CPU   ARM A9 CPU
Application Processor Unit (APU)
DMA channel
Central Interconnect

clock

AXI ACP slave port   GP AXI master port   GP AXI slave port

I/O Peripherals   COM

DDR2/3 controller   1 GB

HP AXI slave Port

**links**

programmable logic (PL)

BRAM   SB   logic   SB   logic   SB   logic   SB   ...   SB   DSP   SB   logic   SB

BRAM   SB   logic   SB   logic   SB   logic   SB   ...   DSP   SB   logic   SB

BRAM   SB   SB   SB   SB   DSP   SB

DSP

PCle   Host

12x   GTX

→ „Up"
→ „Down"
→ „North"
→ „West"
→ „East"
→ „South"

**16 serial transceivers**,
- up to 12.5 Gb/s data rate (GTX)

**545 Block RAMs** (à 36 kb) ,
- total 19.2 Mbit
- configurable as single/dual port SRAM
- up to 72 bit data width

**900 DSP slices**
- 18x25 bit signed multiply
- 48-bit adder/accumulator
**~ 180 single float Adder/Multiplier**

# NODE SCHEMATICS, MODELS, AND DESIGN FLOW



**Spike Receive & Connections**

**Solving Neuro-Synaptic Equations & Spike Detection**

CB → SS → state memory ← ES

D | W | N | T

$V_m$ | $I_e$ | $I_i$

RTR

ODE

EXT / seed

$\delta_p$

$\delta_q$

FIFO 256 words

FIFO 512 words | AER-in

AER-out

AXI-W

AXI-R

1GB DDR

PD — PE

**(64 Bit)** input packet | output packet **(64 Bit)**

node router

„East"  „North"  „Up"  „Down" „South"  „West"
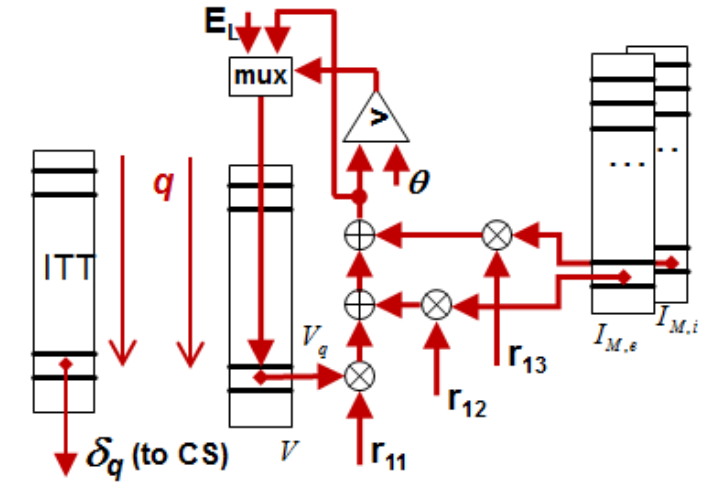
**Simulation Node Schematics**

- **256-Neuron-Fascicle per Node**

- **Various Point-Neuron Models**
  - **- LIF, MAT2, Izhikevich, AdEx**

- **Several Synapse Models**
  - - CUBA- and COBA-based,
  - - Exponential Decay-, Alpha-, Beta-Shaped

- **Several ODE System Solvers**
  - - Exact Exponential
  - - Runge-Kutta
  - - Parker-Sochacki

- **Network Generation „on the fly" during simulation run time based on highly efficient PRNGs („Procedural Connectivity")**

- **Node-Synchronization by Barrier-Messages**
  - - avoid spike-loss

- **Single-Float Precision Arithmetics**

- **High-Level-Synthesis Design Flow**

# ODE SOLVER

- membrane of a point-neuron (LIF)

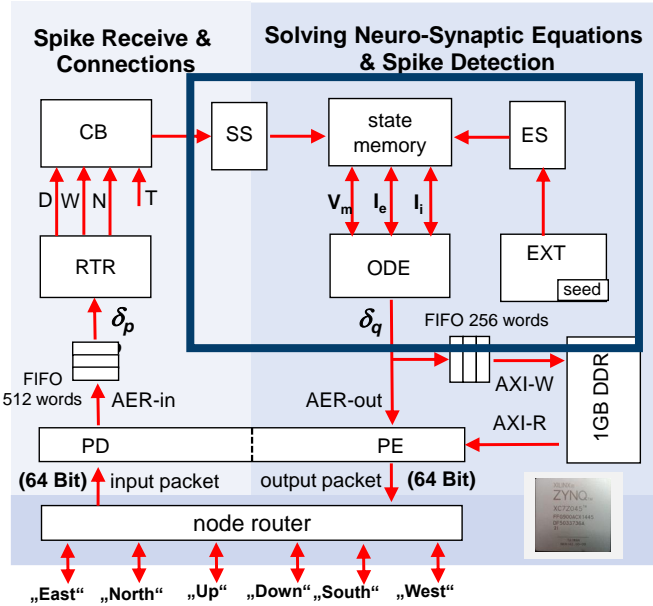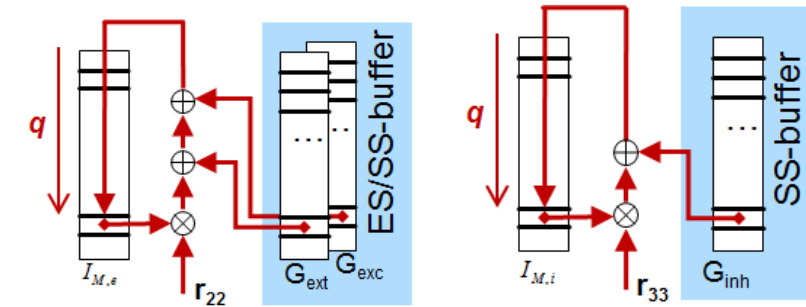$$\tau_m \cdot \frac{dV_q}{dt} = -(V_q - E_L) + R_M \cdot I_{M,q}$$

synthesized datapaths, pipelined

- lumped CUBA-synapses

  - exponential decay

$$\tau_{s,x} \cdot \frac{dI_{M,q,x}}{dt} = -I_{M,q,x} + I_{S,q,x} \quad , \quad x \in \{e,i\}$$

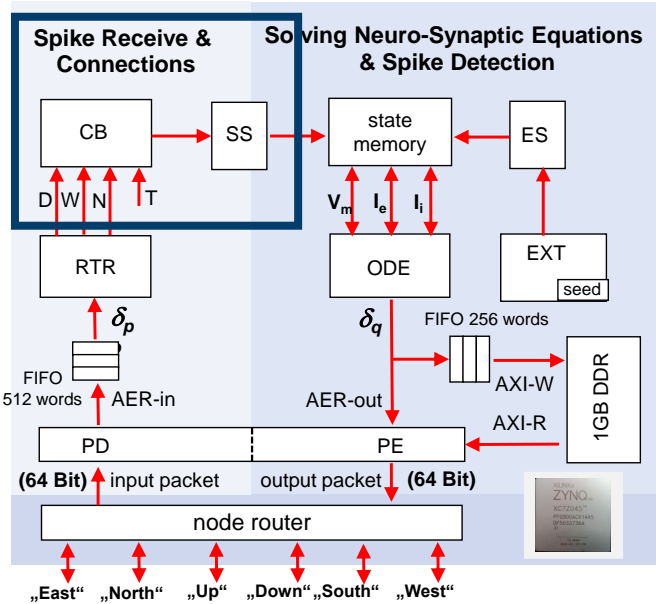- State-Vector update

  - exact exponential for linear ODEs

$$\begin{bmatrix} V_q \\ I_e \\ I_i \end{bmatrix}(t+h) \leftarrow \begin{bmatrix} r_{11} & r_{21} & r_{31} \\ 0 & r_{22} & 0 \\ 0 & 0 & r_{33} \end{bmatrix} \cdot \begin{bmatrix} V_q \\ I_e \\ I_i \end{bmatrix}(t)$$
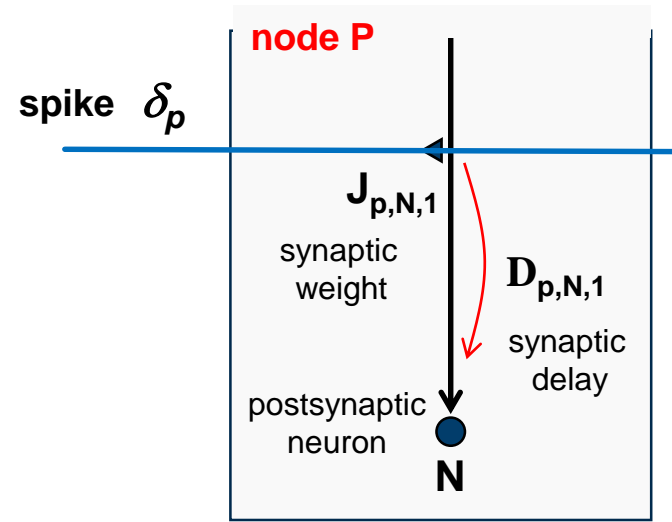
update equation, 3 state variables per neuron

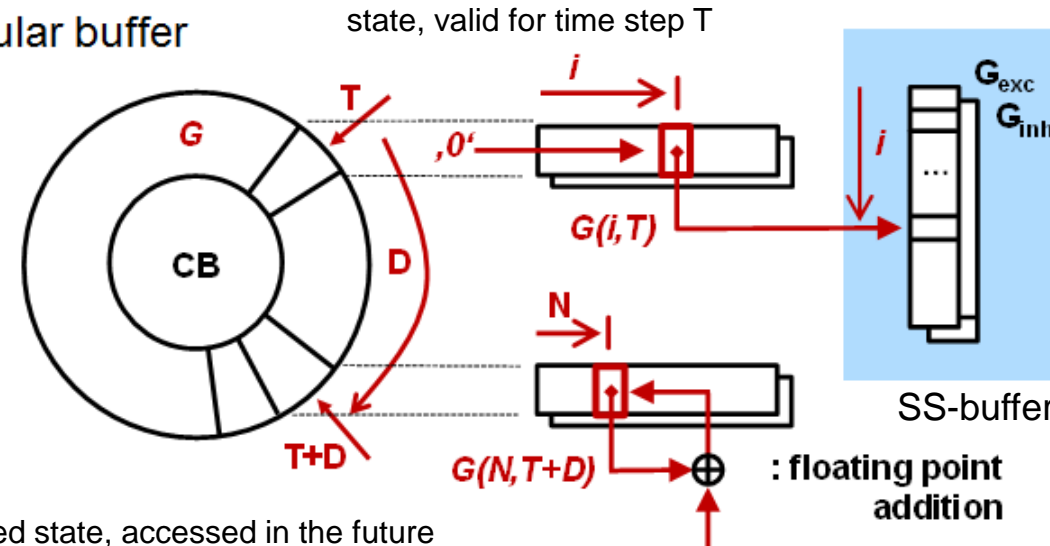# LUMPED SYNAPSES

- network schematic

- synaptic equation

**node P**

spike $\delta_p$

$J_{p,N,1}$
synaptic weight

$D_{p,N,1}$
synaptic delay

postsynaptic neuron

**N**

**synaptic multiplicity (multapses)**

**pre-synaptic spike train**

$$I_{S,q,x}(T) = \sum_{p \in B_x} \sum_{m=1}^{M_{q,p}} J_{q,p,m} \cdot S_p(T - D_{q,p,m})$$

**synaptic weight**

**synaptic delay**

$$S_p(t) = \sum_i \delta(t - t_{p.i})$$

- circular buffer

state, valid for time step T

$G$

**CB**

$D$

$T$

$i$

,0'

$G(i,T)$

$G_{exc}$

$G_{inh}$

$i$

+ addressing by ,time' **T**

and ,neuron index' **N**

SS-buffer

$N$

$T+D$

$G(N,T+D)$

$\oplus$ : floating point addition

delayed state, accessed in the future

Spike Receive & Connections

Solving Neuro-Synaptic Equations & Spike Detection

CB — SS — state memory — ES

D W N T

$V_m$  $I_e$  $I_i$

RTR

EXT
seed

$\delta_p$

$\delta_q$

FIFO 256 words

ODE

FIFO 512 words — AER-in

AER-out

AXI-W

AXI-R

1GB DDR

PD

PE

(64 Bit) input packet   output packet (64 Bit)

node router

„East"  „North"  „Up"  „Down" „South"  „West"

# PROCEDURAL CONNECTIVITY

address    memory

- **Parameters of a synapse C**

- **example network**

- **(naive) memory layout on the receiving node**

$C_{k,j}$   $C_{p,j}$   $C_{q,j}$

presynaptic neuron: j

postsynaptic neurons: k   p   q

t+D

C   N

$A_j$

higher addresses

$D_{k,j}$
$W_{k,j}$
$k$   $\}$ $C_{k,j}$
$D_{p,j}$
$W_{p,j}$
$p$   $\}$ $C_{p,j}$
$D_{q,j}$
$W_{q,j}$
$q$   $\}$ $C_{q,j}$

C
- D    ax. + dend. delay
- W    weight
- N    postsynaptic neuron

JÜLICH
Forschungszentrum

# PROCEDURAL CONNECTIVITY



HP slave Port

$A_j \longrightarrow$ X $\longrightarrow$ DDR2/3 CTRL $\longrightarrow$ A SDRAM (1GB)

$C_j \longleftarrow$ X $\longleftarrow$ DDR2/3 CTRL $\longleftarrow$ D

$$\boxed{\begin{array}{c} D_{p,j} \\ W_{p,j} \\ p \end{array}}$$

$A_j$

$t_{access}$ (~ 200 ns)

**~30 cycles @150 MHz**

$C_j$

t (time)

- access latency

- minimal time from „address valid" to „first data out"

# NETWORK GENERATION

**Microcircuit:** connections $C$ are defined
by pseudo random numbers



D,W: normal-distribution

N: binomial-distribution

weight distribution



excitatory weights

L4E→L23E

Inhib.

delay distribution



Inhib.

Excit.

- In NEST: all synapse parameters are drawn offline using pseudo-random-number generators (PRNGs)

- **Alternative approach**: implement these random number generators on-the-chip

  - re-generate the synaptic parameters „on-the-fly", when needed

  - initial seeds define the deterministic sequence of random numbers

seed($D_i$) → PRNG → D

# PROCEDURAL CONNECTIVITY

**seed-memory (on-chip BRAM)**

address $A_j$

seed$(D_j)$
seed$(W_j)$
seed$(N_j)$
$L_j$

seed$(D_j)$ → PRNG → D
$L_i$ X

seed$(W_j)$ → PRNG → W
$L_i$ X

seed$(N_j)$ → PRNG → N
$L_i$ X

**Pseudo-Random-Number Generators: create sequences in parallel**

clock · AXI ACP slave port · GP AXI m

BRAM · SB · logic · SB · logic · SB · log
BRAM · SB · logic · SB · logic · SB · log
BRAM · SB · SB · SB

XILINX ZYNQ XC7Z045 FFG900ACX1445 2I

**Timing**

$A_j$

1 cycle @150 MHz

seed

2 cycle @150 MHz

| $D_{k,j}$ | $D_{p,j}$ | $D_{q,j}$ |
| $W_{k,j}$ | $W_{p,j}$ | $W_{q,j}$ |
| k | p | q |

1 cycle @150 MHz

t (time)

# PROCEDURAL CONNECTIVITY

- seed-lookup tables



AER-Address of presynaptic neuron $\delta_p$

Seed-Tables

- 96 seed-bits per presynaptic neuron

- compression factor ~ 6.5

- full connectome can be stored in the local BRAM

- RTR data path



- Walker's table-based PRNG

- 64 sampling points per CDF/PDF

Walker, A. J., "An Efficient Method for Generating Discrete Random Variables with General Distributions". ACM Transactions on Mathematical Software. **3** (3): 253–256, 1977, doi:10.1145/355744.355749

- RTR control loop



- latency: 3 cycles @ 150 MHz

- fully pipelined

# Gate and Memory Breakdown (LIF, CUBA)

- Overall FPGA resources



**Look-up tables**
$N_{LUT}$ = 218.600

51 % · 37 % · 11 %

**Flip-Flops**
$N_{FF}$ = 437.200

73 % · 19 % · 8 %

**Block-RAM**
$N_{BRAM}$ = 19.2 Mbit

62 % · 26 % · 12 %

Legend:
- CnNN
- Node-Router
- free (unused)

- Circuit resources for CnNN



**Look-up tables**
$N_{LUT}$ = 23.860

4 % · 35 % · 28 % · 8 % · 16 % · 9 %

**Flip-Flops**
$N_{FF}$ = 32.662

2 % · 22 % · 9 % · 9 % · 9 % · 49 %

**Block-RAM**
$N_{BRAM}$ = 12 Mbit

9 % · 12 % · 79 %

Legend:
- ODE
- Posson Stimulus
- Network Generation
- Control & Setup
- Communication Layer
- other

- <span style="color:red">SNNs are memory-dominated</span>

# LATENCY-BREAKDOWN : THE MICRO-CIRCUIT



$N_h$ = 2905 cycles

$N_h$ = 3690 cycles

$N_h$ = 6452 cycles

fast case

average
($\mu$ = 24.6 $\mu$s)

worst case

■ : **Communication**    ■ : **Poisson-Stimulus**

■ : **Network Generation**    ■ : **ODE**

- procedural connectivity: 20 % speedup over ext. DRAM

- Speedup  X 4.06  over BRT

# TOWARDS REALISTIC NEURON MODELS

|  | The 'common case' | Computational Complexity |
|---|---|---|
| **present** | • **Linear** point-neuron dynamics <br><br> • **No dendritic tree** (galvanic interconnect) | • ODE integration by **exponential Euler's method** <br><br> • **High degree of parallelism** |
| **prospective** | • **Nonlinear** neuron dynamics (e.g. Izhikevich, AdEx, HH) <br><br> • **compartmental dendritic tree**, model of 'the dendritic computational toolkit' *) <br><br> *) M.London, M.Häusser, Annu.Rev.Neurosci. 2005.28:503-532 | • ODE integration by **advanced numeric methods** <br><br> • **Limited degree of parallelism** and significant critical paths in arithmetics <br><br> • particular strategies for **stiff problems** |

Galvanic interconnect

dendrite

point-neuron

RCG-T-section

$g_{link,i+1}$   $V_{i+1}$

$E_{ex}$   $g_{ex,i}$   $g_{m,i}$

$c_{m,i}$

$E_{inh}$   $g_{inh,i}$

$g_{link,i}$   $V_i$

$V_{i-1}$

COBA compartment

# TOWARDS REALISTIC NEURON MODELS

**[Hodkin, Huxley, 1952]**

$$C \cdot \frac{d}{dt}V(t) = -\sum_k I_k(t) + I(t)$$

$$\sum_k I_k(t) = g_{Na} \cdot m^3 \cdot h \cdot [V(t) - E_{Na}] +$$

$$+ g_K \cdot n^4 \cdot [V(t) - E_K] + g_L \cdot [V(t) - E_L]$$

$$\frac{d}{dt}m(t) = \alpha_m(V) \cdot [1 - m(t)] - \beta_m(V) \cdot m(t)$$

$$\frac{d}{dt}n(t) = \alpha_n(V) \cdot [1 - n(t)] - \beta_n(V) \cdot n(t)$$

$$\frac{d}{dt}h(t) = \alpha_h(V) \cdot [1 - h(t)] - \beta_h(V) \cdot h(t)$$

$$\alpha_m(V) = \frac{2.5 - 0.1 \cdot V(t)/mV}{[\exp(2.5 - 0.1 \cdot V(t)/mV) - 1]}$$

$$\beta_m(V) = 4 \cdot \exp(-V(t)/18mV)$$

**IAF-or-Burst (IFB) [Smith et al., 2000]**

$$C \cdot \frac{d}{dt}V(t) = g_L \cdot (V_L - V(t)) + g_T \cdot h \cdot \chi(V(t) - V_h) + I(t)$$

$$\frac{d}{dt}h(t) = \begin{cases} -\dfrac{h(t)}{\tau^-} & V \geq V_h \\ \dfrac{1 - h(t)}{\tau^+} & V < V_h \end{cases}$$

$$V(t^+) = V_r \text{ if } V(t^-) = V_{th}$$

**Adaptive Exponential IAF
[Gerstner, Brette. 2009]**

$$C \cdot \frac{d}{dt}V(t) = g_L \cdot (V(t) - V_R) + g_L \cdot \Delta_T \cdot e^{\frac{V(t) - \vartheta_{th}}{\Delta_T}} - w(t) + I(t)$$

$$\tau \cdot \frac{d}{dt}w(t) = a \cdot (V(t) - E_L) - w(t)$$

$$V(t^+) = V_r \text{ if } V(t^-) = V_{th}$$

**Quadratic IAF
[Latham et al., 2000]**

$$C \cdot \frac{d}{dt}V(t) = a_0 \cdot (V(t) - V_R) \cdot (V(t) - V_C) + I(t)$$

$$V(t^+) = V_r \text{ if } V(t^-) = V_{th}$$

**Izhikevich
[Izhikevich, 2003]**

$$\frac{d}{dt}u(t) = 0.04 \cdot u^2(t) + 5 \cdot u(t) + 140 - w(t) + I(t)$$

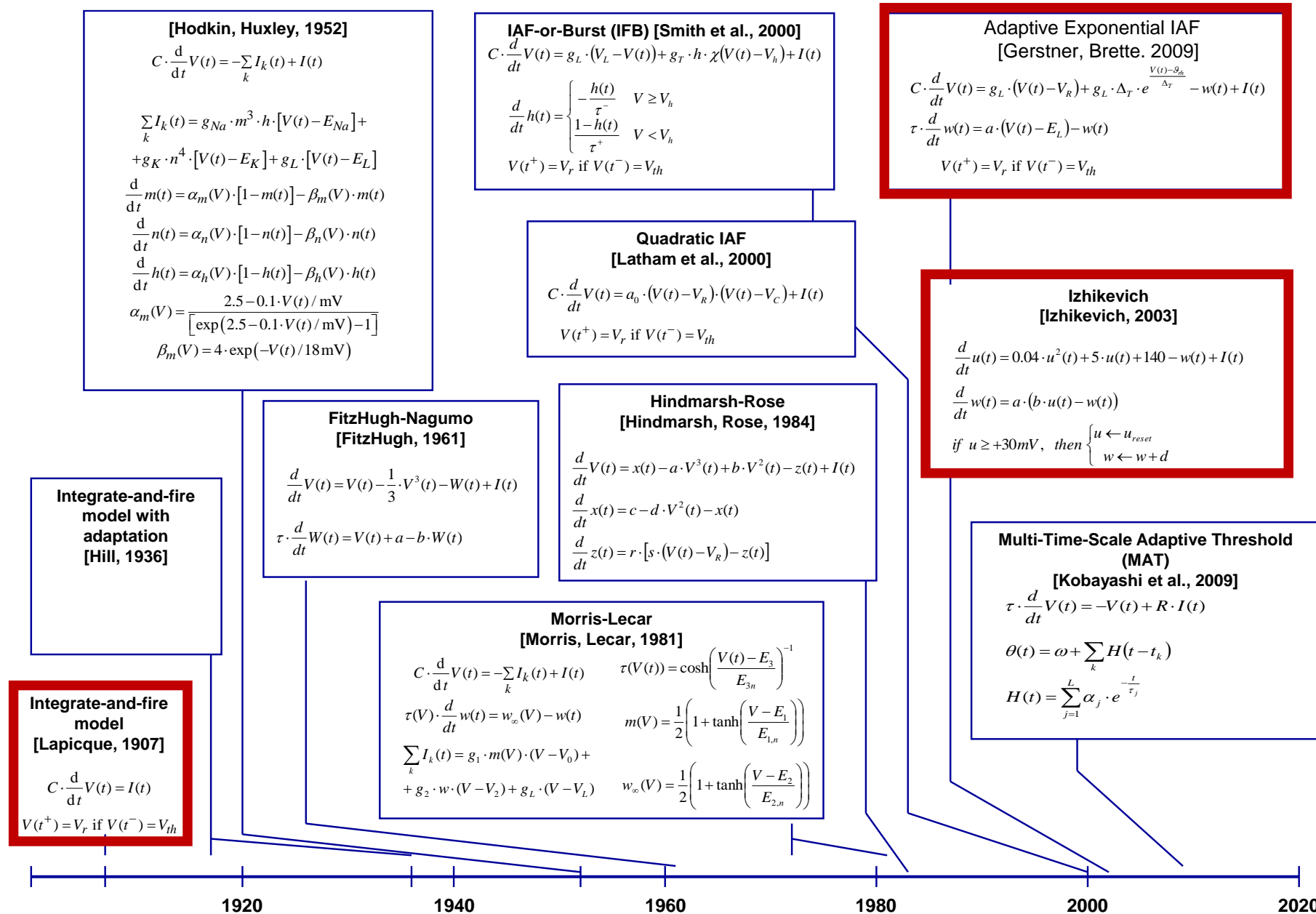$$\frac{d}{dt}w(t) = a \cdot (b \cdot u(t) - w(t))$$

$$\text{if } u \geq +30mV, \text{ then } \begin{cases} u \leftarrow u_{reset} \\ w \leftarrow w + d \end{cases}$$

**FitzHugh-Nagumo
[FitzHugh, 1961]**

$$\frac{d}{dt}V(t) = V(t) - \frac{1}{3} \cdot V^3(t) - W(t) + I(t)$$

$$\tau \cdot \frac{d}{dt}W(t) = V(t) + a - b \cdot W(t)$$

**Hindmarsh-Rose
[Hindmarsh, Rose, 1984]**

$$\frac{d}{dt}V(t) = x(t) - a \cdot V^3(t) + b \cdot V^2(t) - z(t) + I(t)$$

$$\frac{d}{dt}x(t) = c - d \cdot V^2(t) - x(t)$$

$$\frac{d}{dt}z(t) = r \cdot [s \cdot (V(t) - V_R) - z(t)]$$

**Integrate-and-fire
model with
adaptation
[Hill, 1936]**

**Multi-Time-Scale Adaptive Threshold
(MAT)
[Kobayashi et al., 2009]**

$$\tau \cdot \frac{d}{dt}V(t) = -V(t) + R \cdot I(t)$$

$$\theta(t) = \omega + \sum_k H(t - t_k)$$

$$H(t) = \sum_{j=1}^L \alpha_j \cdot e^{-\frac{t}{\tau_j}}$$

**Morris-Lecar
[Morris, Lecar, 1981]**

$$C \cdot \frac{d}{dt}V(t) = -\sum_k I_k(t) + I(t)$$

$$\tau(V) \cdot \frac{d}{dt}w(t) = w_\infty(V) - w(t)$$

$$\sum_k I_k(t) = g_1 \cdot m(V) \cdot (V - V_0) +$$

$$+ g_2 \cdot w \cdot (V - V_2) + g_L \cdot (V - V_L)$$

$$\tau(V(t)) = \cosh\left(\frac{V(t) - E_3}{E_{3n}}\right)^{-1}$$

$$m(V) = \frac{1}{2}\left(1 + \tanh\left(\frac{V - E_1}{E_{1,n}}\right)\right)$$

$$w_\infty(V) = \frac{1}{2}\left(1 + \tanh\left(\frac{V - E_2}{E_{2,n}}\right)\right)$$

**Integrate-and-fire
model
[Lapicque, 1907]**

$$C \cdot \frac{d}{dt}V(t) = I(t)$$

$$V(t^+) = V_r \text{ if } V(t^-) = V_{th}$$

1920    1940    1960    1980    2000    2020

ACA

**JÜLICH**
Forschungszentrum

# ODE SOLVERS

**Exact Integration (applies for linear ODEs only)**

$$V(t+h) = e^{Ah} \cdot V(t)$$

**S-step Runge-Kutta Method   (RK-S)**

$$k_j = f\left( t_n + h \cdot c_j \; ; \; V(t) + h \cdot \sum_{l=1}^{s} a_{jl} \cdot k_l \right), j = 1,...,s$$

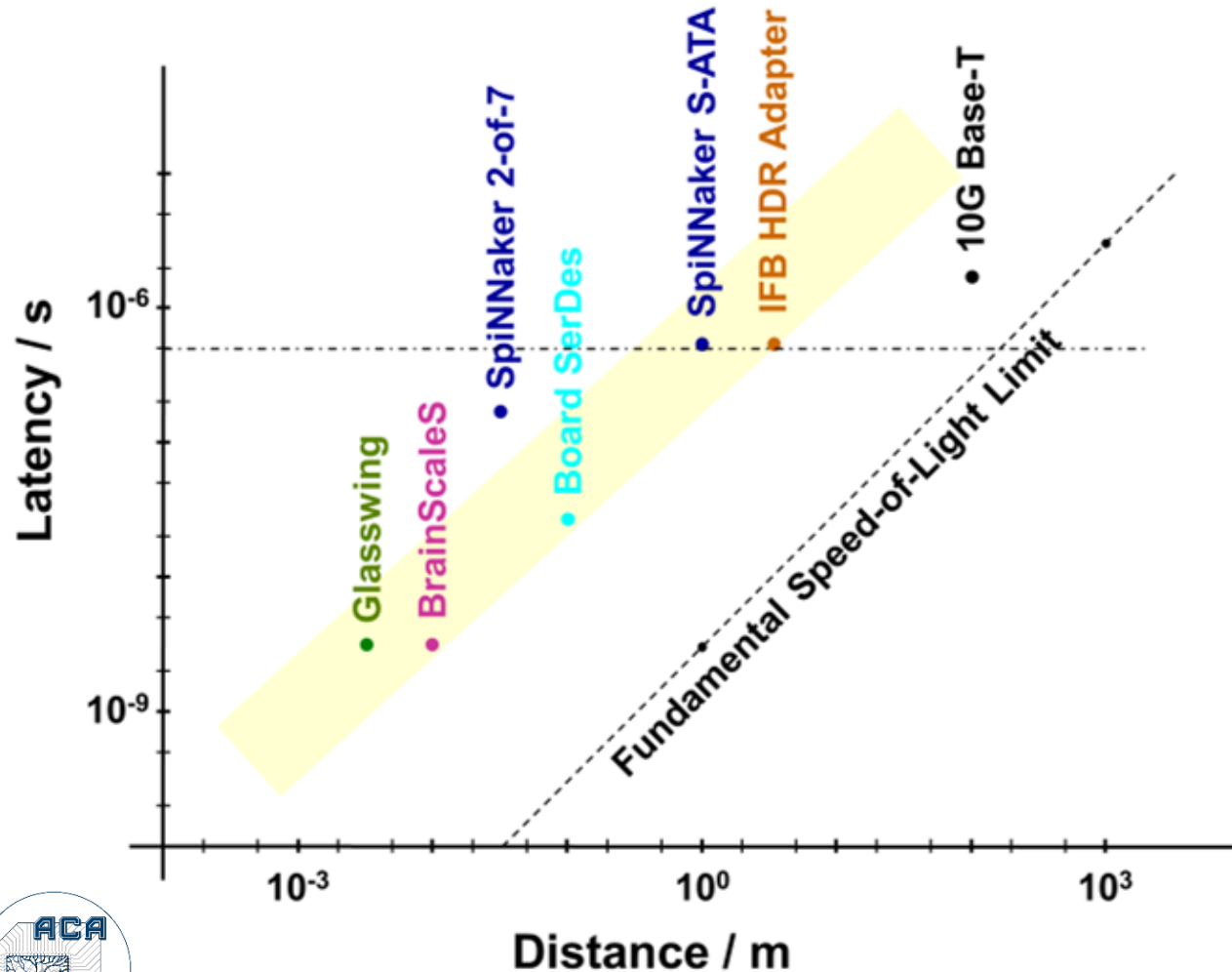$$V(t+h) = V(t) + h \cdot \sum_{j=1}^{s} b_j \cdot k_j$$

**S-step Parker-Sochacki Method (PS-S)**

$$V_i(t+h) = V_i(t) + \frac{V_i'(t)}{1!} \cdot h + \frac{V_i''(t)}{2!} \cdot h^2 + \frac{V_i'''(t)}{3!} \cdot h^3 + ... + \frac{V_i^{(s)}(t)}{s!} \cdot h^s$$

⇒ **Speed-up not limited by arithmetics**

# Trends in State-of-the-Art Communication Standards
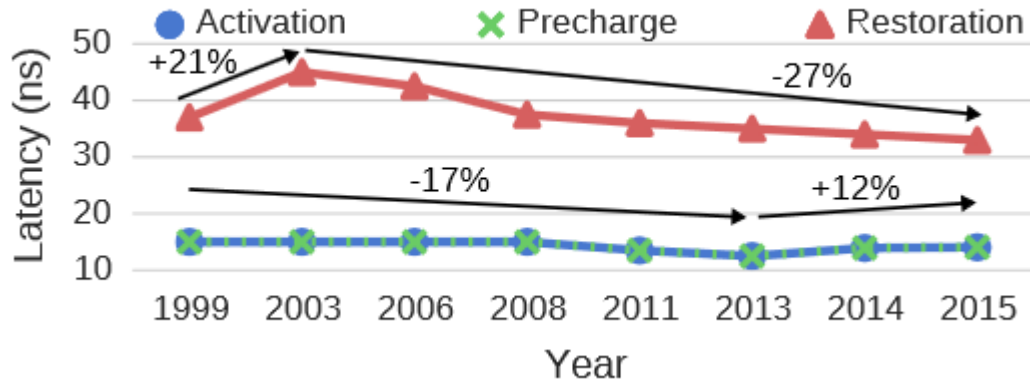
**1-Hub latency**



- Latency grows with the distance of compute-nodes

- Avoid large latencies by

  - few hubs between nodes

  - dense (3D-stacked) circuit integration

- Conceptually, highly integrated components with short physical distance could be required

# Trends in State-of-the-Art DRAM Performance

- Performance data from 54 different commercially available SDRAMs (2022)

- Latency Trends



K.K.Change et al., „Understanding Latency Variation in modern DRAM Chips: Experimental Characterization, Analysis, and Optimization", doi:10.1145/2896377.2901453 (2016)
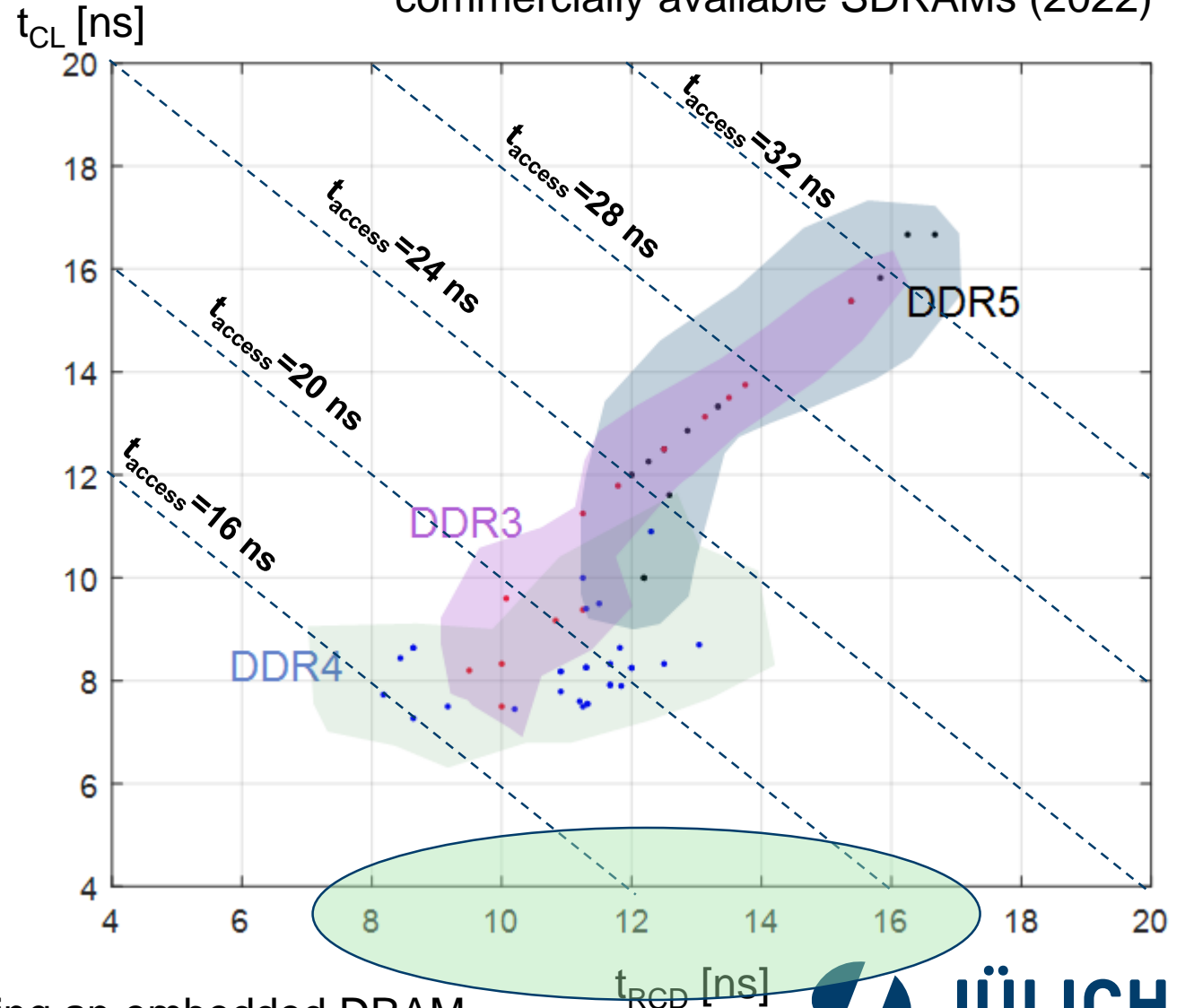
- No significant trend in access latency reduction

- Neuromorphic Computing:

  - **Plasticity** will incur significant traffic to the memory system

  - latency issue could be solved e.g. by using an embedded DRAM technology in a dedicated accelerator circuit
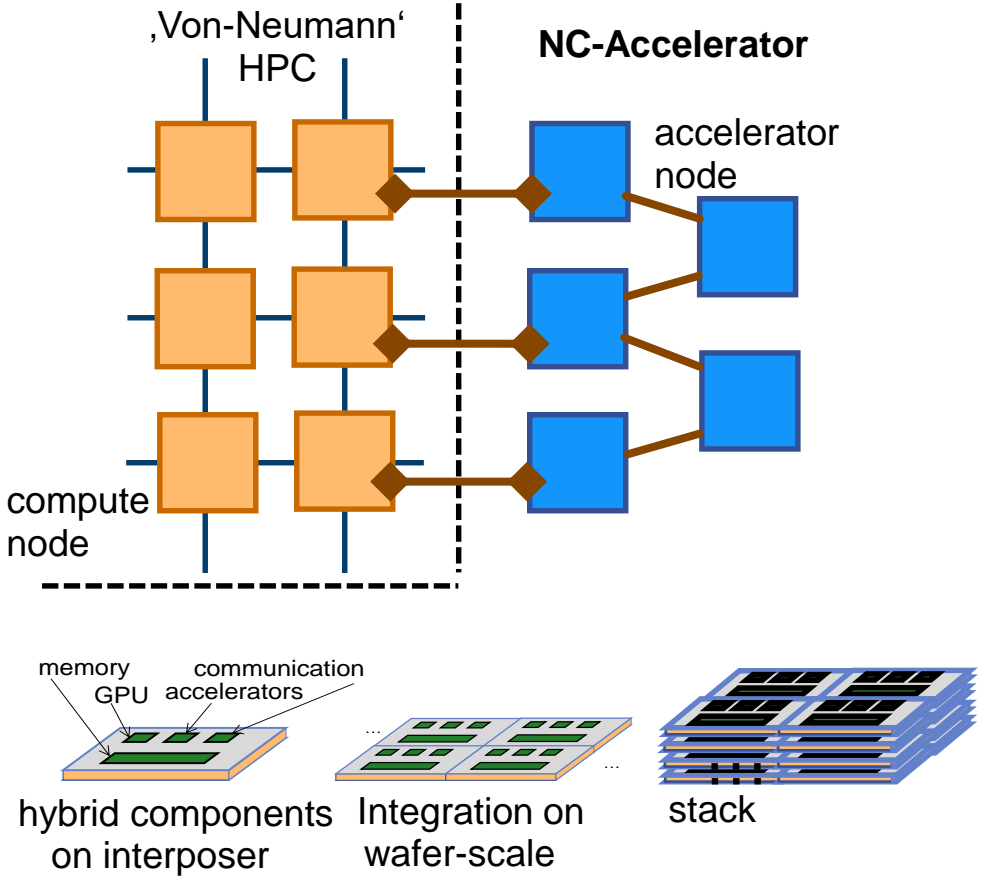
# CONCLUSION

- Components in conventional HPC systems are optimized for both:

  **large data packets** and **large communication & memory bandwidth**

  → results in **large latencies** for …. memory accesses and packet transmission


- Requirements for a future accelerated neuromorphic compute platform

  + communication of **data packets** comprising **small size** (spikes)

  + fully random memory access to **small data packets** (e.g. synapse parameters) with very restricted

    locality properties (caches won't work efficiently)

  + **ultra-short latencies for communication and memory access**


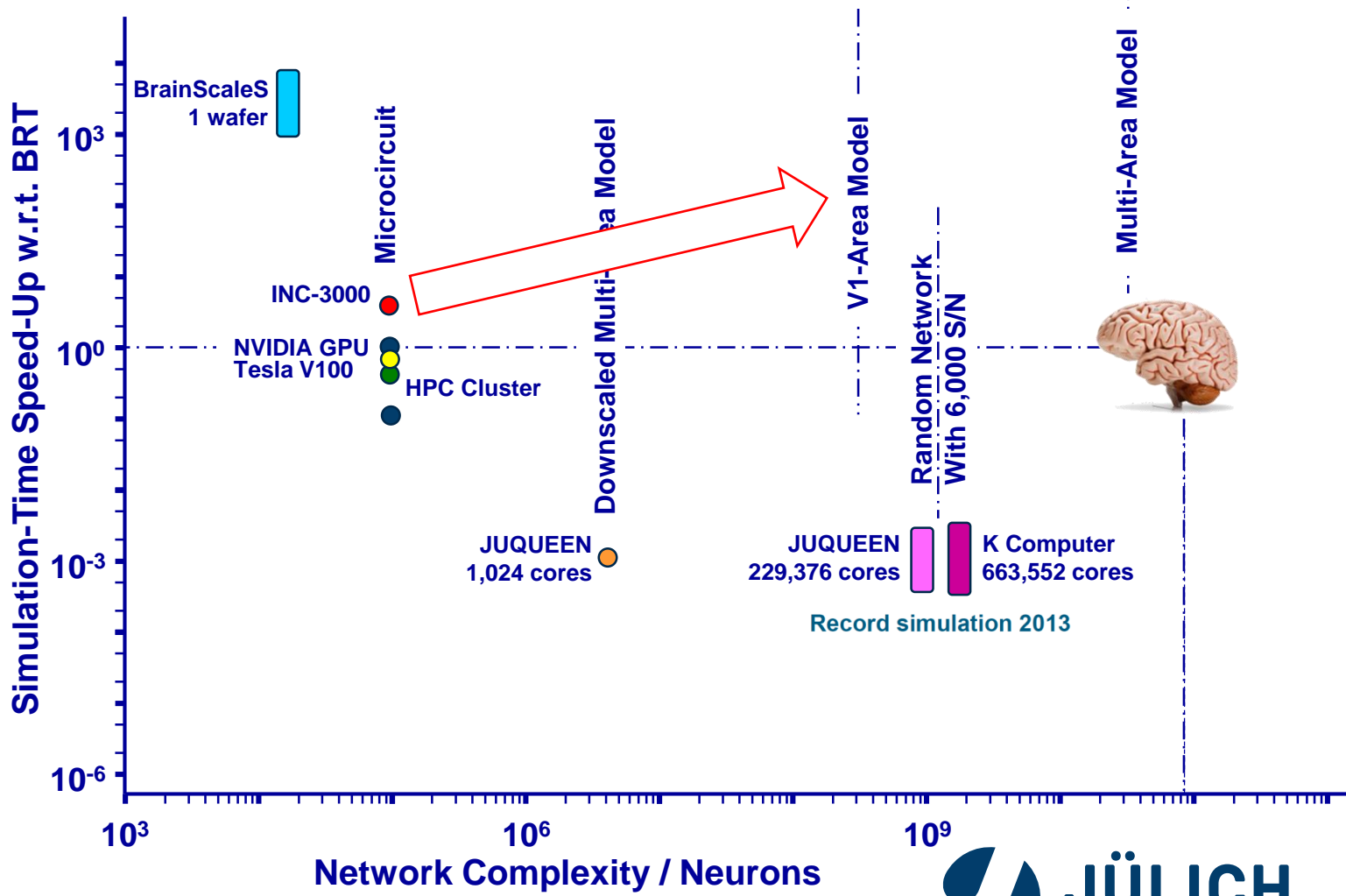The requirements for future NC-Platform and HPC are complementary

- New memory architecture *: near-memory computation with short memory latencies*
- Hierarchical networks with ultra-short communication latency
- *High-package density (2.5D/3D stacked silicon on interposer)*

# THE LONG-TERM GOAL: SUSTAINABLY STAY AHEAD



,Von-Neumann' HPC

NC-Accelerator

accelerator node

compute node

memory GPU   communication accelerators

hybrid components on interposer

Integration on wafer-scale

stack

**State-of-the-Art Simulation Time** (0.1-ms time grid)

Simulation-Time Speed-Up w.r.t. BRT

Network Complexity / Neurons

BrainScaleS 1 wafer

Microcircuit

INC-3000

NVIDIA GPU Tesla V100

HPC Cluster

Downscaled Multi-Area Model

V1-Area Model

Multi-Area Model

Random Network With 6,000 S/N

JUQUEEN 1,024 cores

JUQUEEN 229,376 cores

K Computer 663,552 cores

Record simulation 2013

$10^3$

$10^0$

$10^{-3}$

$10^{-6}$

$10^3$

$10^6$

$10^9$

„Hybrid Neuromorphic-von-Neumann general purpose  computing"

**JÜLICH**
Forschungszentrum