

Cognitronics:

Resource-efficient Architectures for Cognitive Systems

Ulrich Rückert

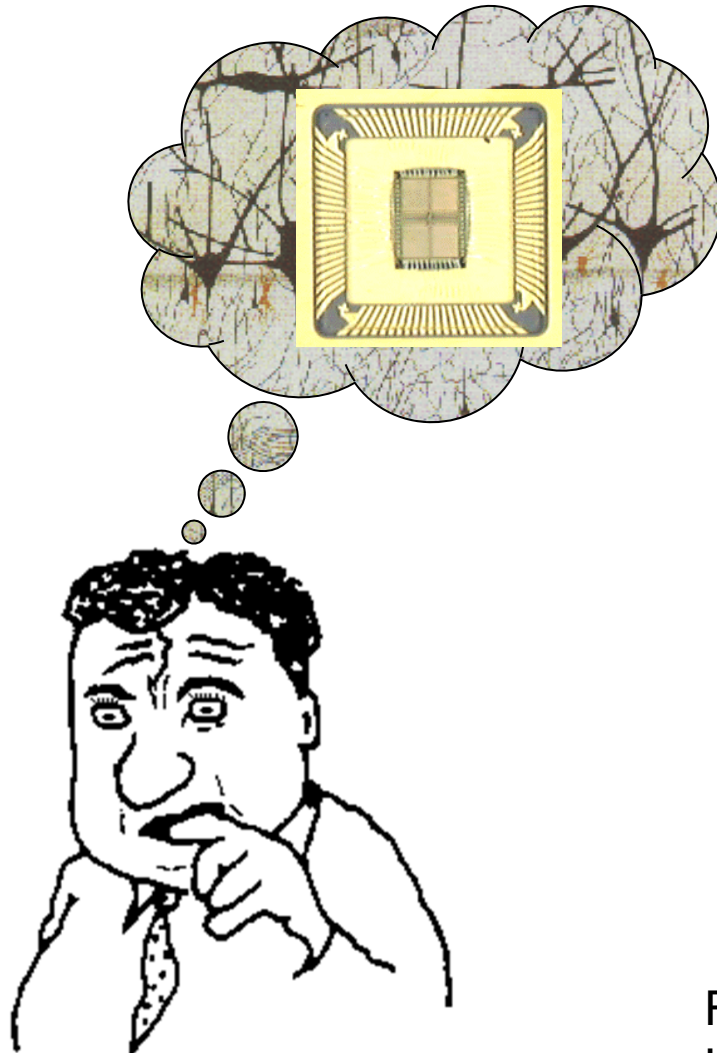
Cognitronics and **S**ensor **S**ystems

14th IWANN, 2017

Cadiz, 14. June 2017

rueckert@cit-ec.uni-bielefeld.de

www.ks.cit-ec.uni-bielefeld.de



Is there a Silicon Way to Neural Networks ?



(1921-1993)

R. Caianiello,

WOPPLOT 1986, Munich

K. Goser / U. Rückert,

IEEE Micro Vol.9, No.6 1989



➤ **Introduction**

Background

➤ **Technology**

Femtoelectronics for ANN

➤ **Architectures**

Design Alternatives

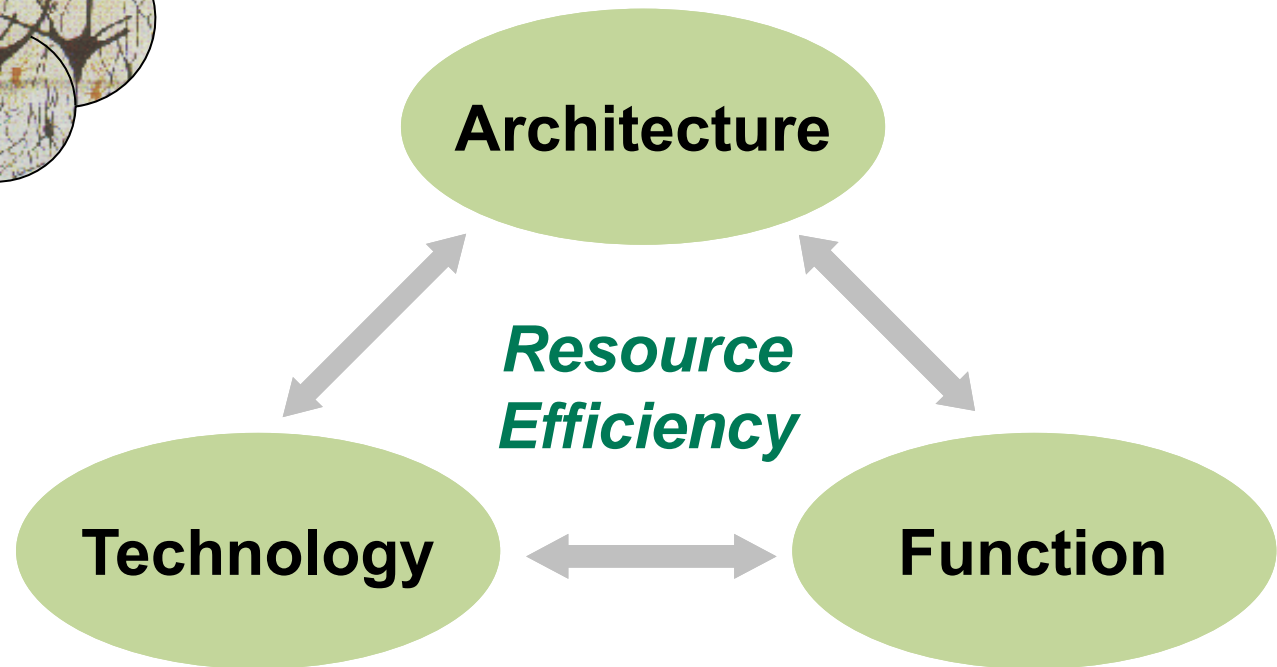
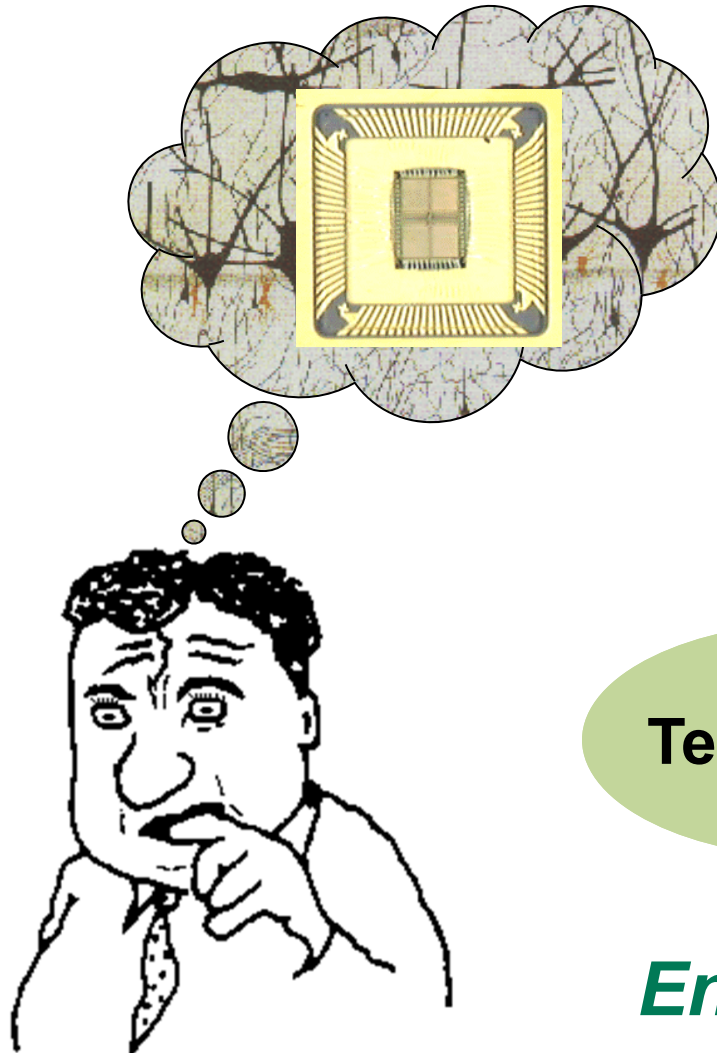
➤ **Applications**

Cognitive Robotics

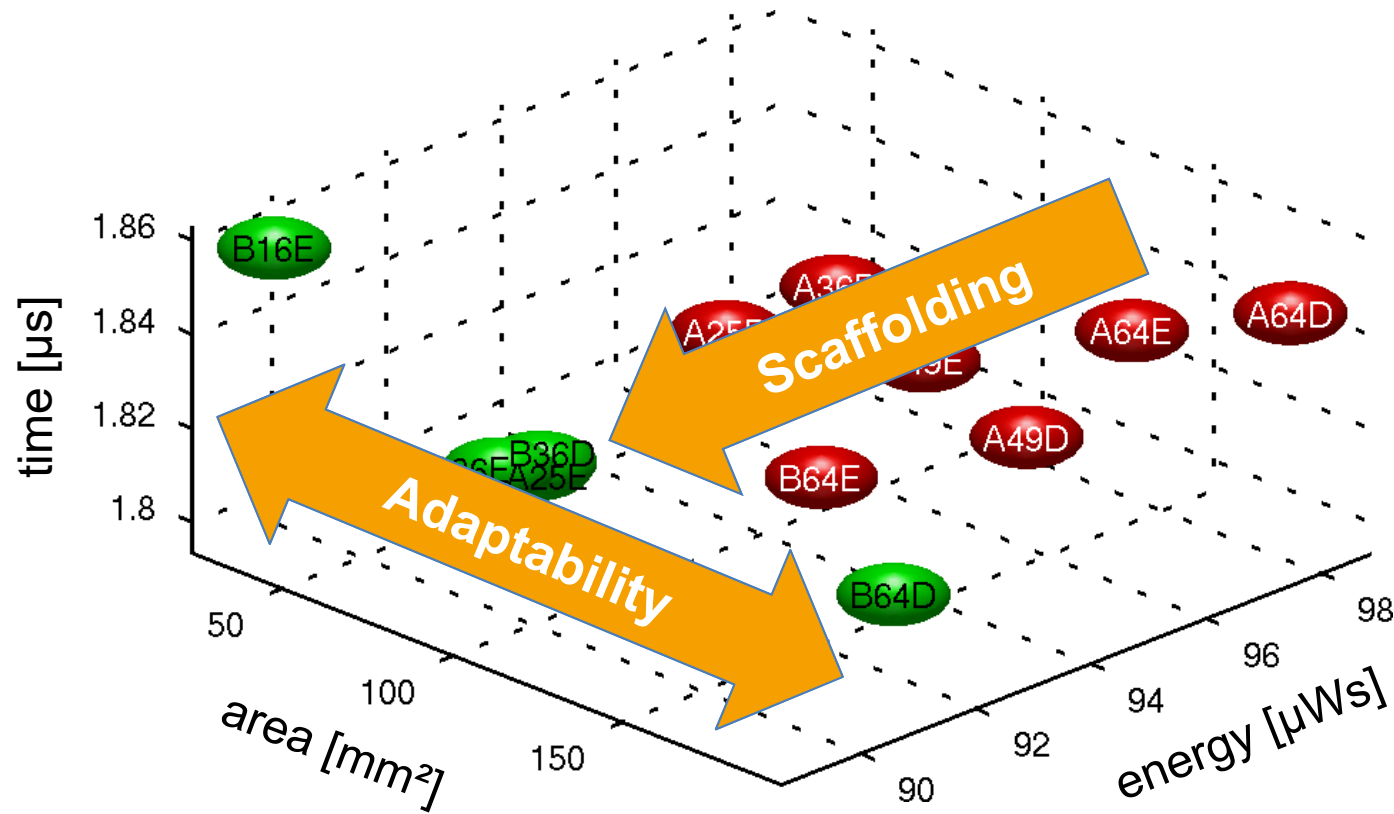
➤ **Discussion**

Questions

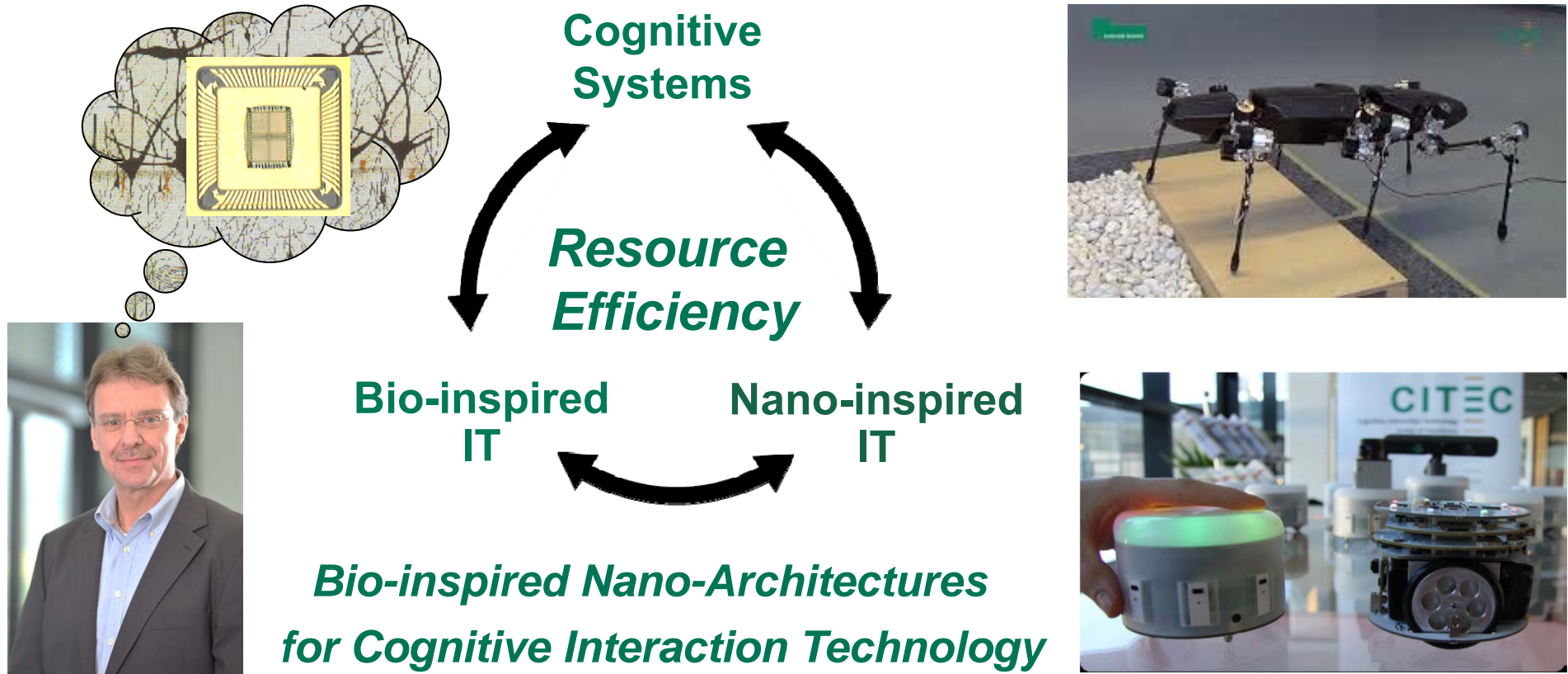
How to use 10^{10} (unreliable) devices/cm² efficiently?



Energy – Mass (Volume) – Time



pareto-optimal design in respect to area, energy, and time



How to generate complex behaviour in the limits of restricted resources?

➤ Introduction

Background



➤ Technology

Femtoelectronics for ANN

➤ Architectures

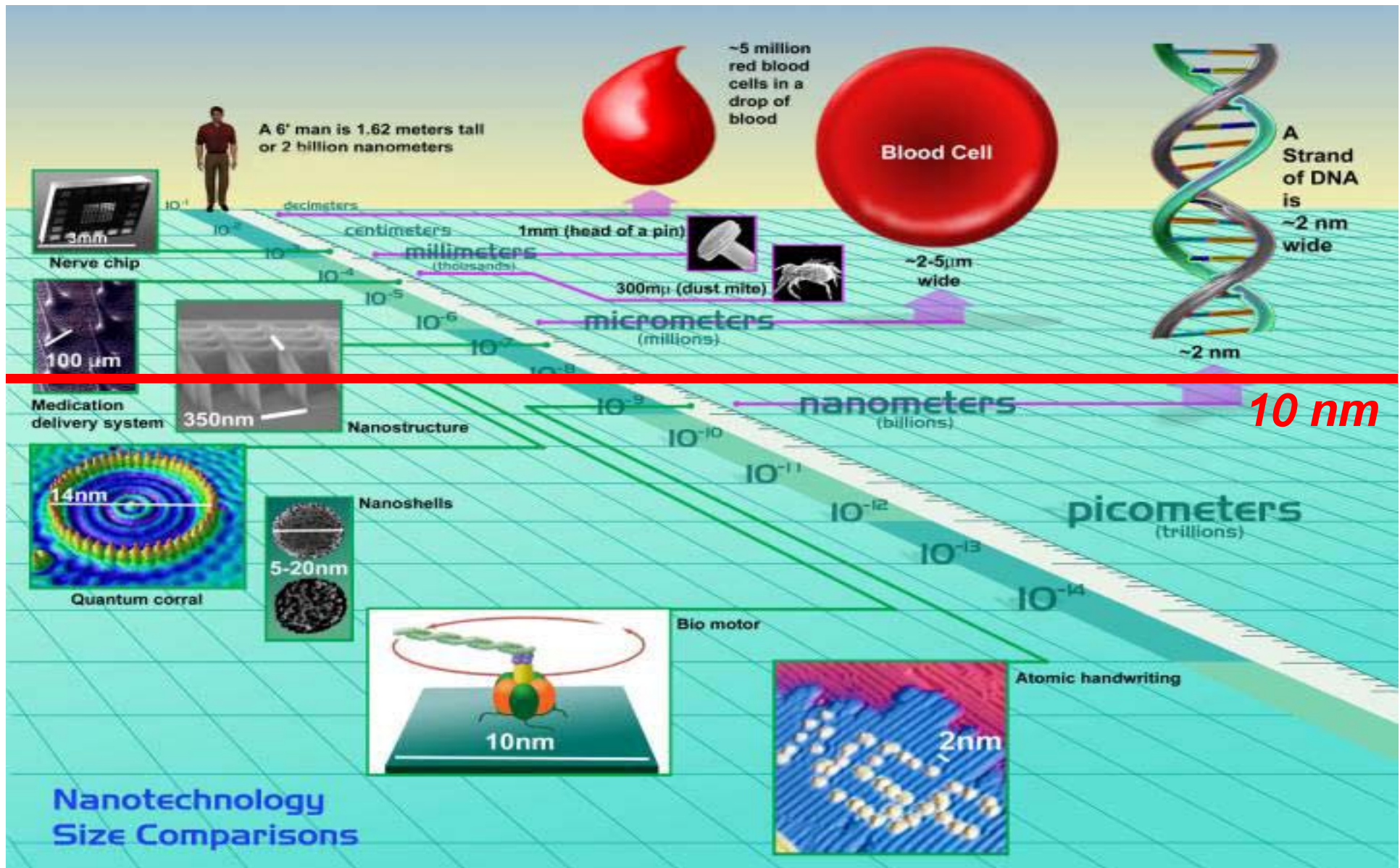
Design Alternatives

➤ Applications

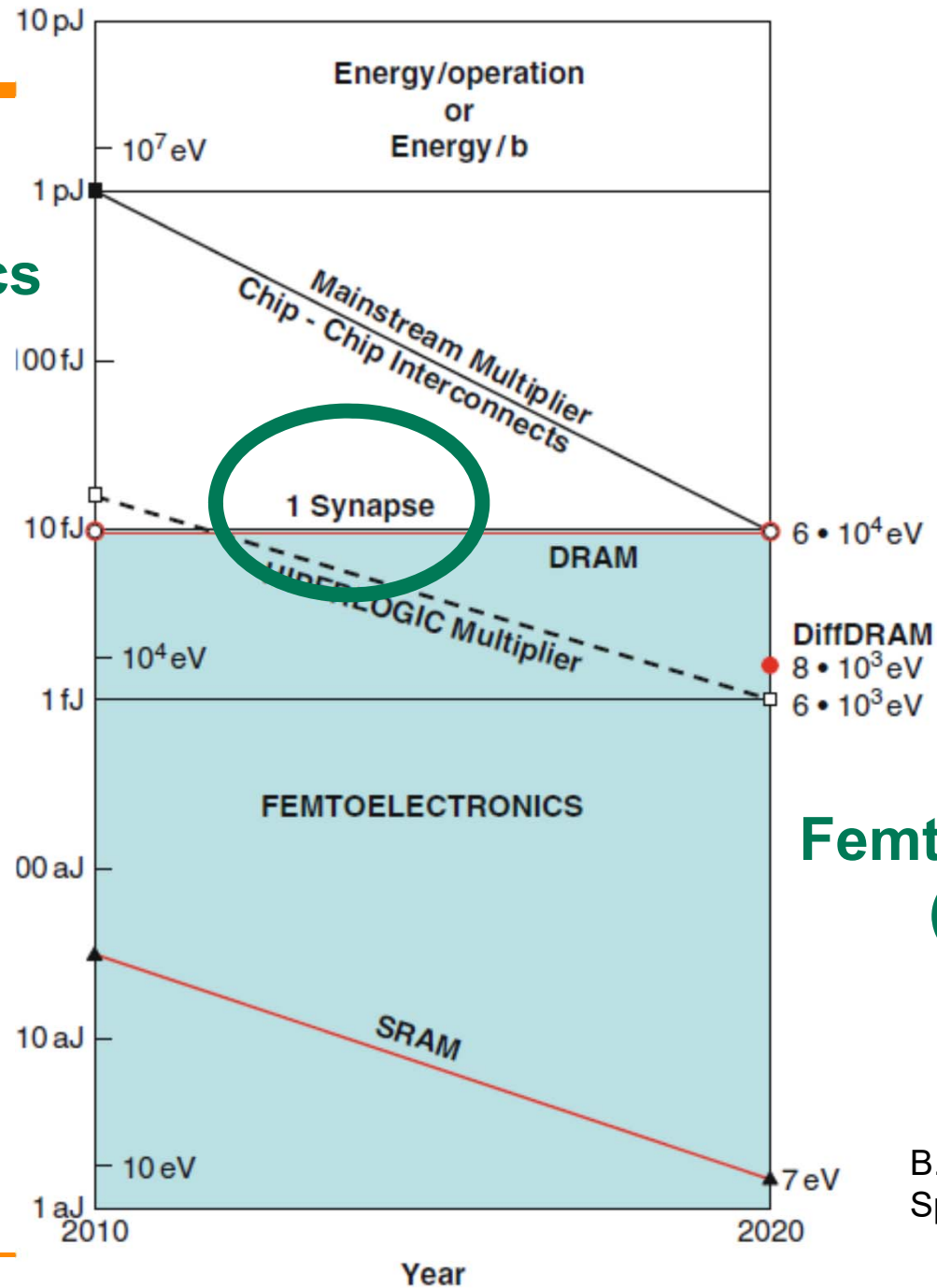
Cognitive Robotics

➤ Discussion

Questions



Nanoelectronics (structure)



Femtoelectronics (energy)

B. Höfflinger: CHIPS 2020
Springer 2012

How to use 10^{10} (unreliable) devices/cm² efficiently?

Need for novel architectures:

- regular & modular structure → reduced design complexity
- fault-tolerance, redundancy → higher yield
- parallelism → reduced power consumption!
- scalability → faster and simpler mapping to next generation technologies

Neural architectures are very good candidates!

Technology Push!

(1990)

➤ Introduction

Background

➤ Technology

Femtoelectronics for ANN



➤ Architectures

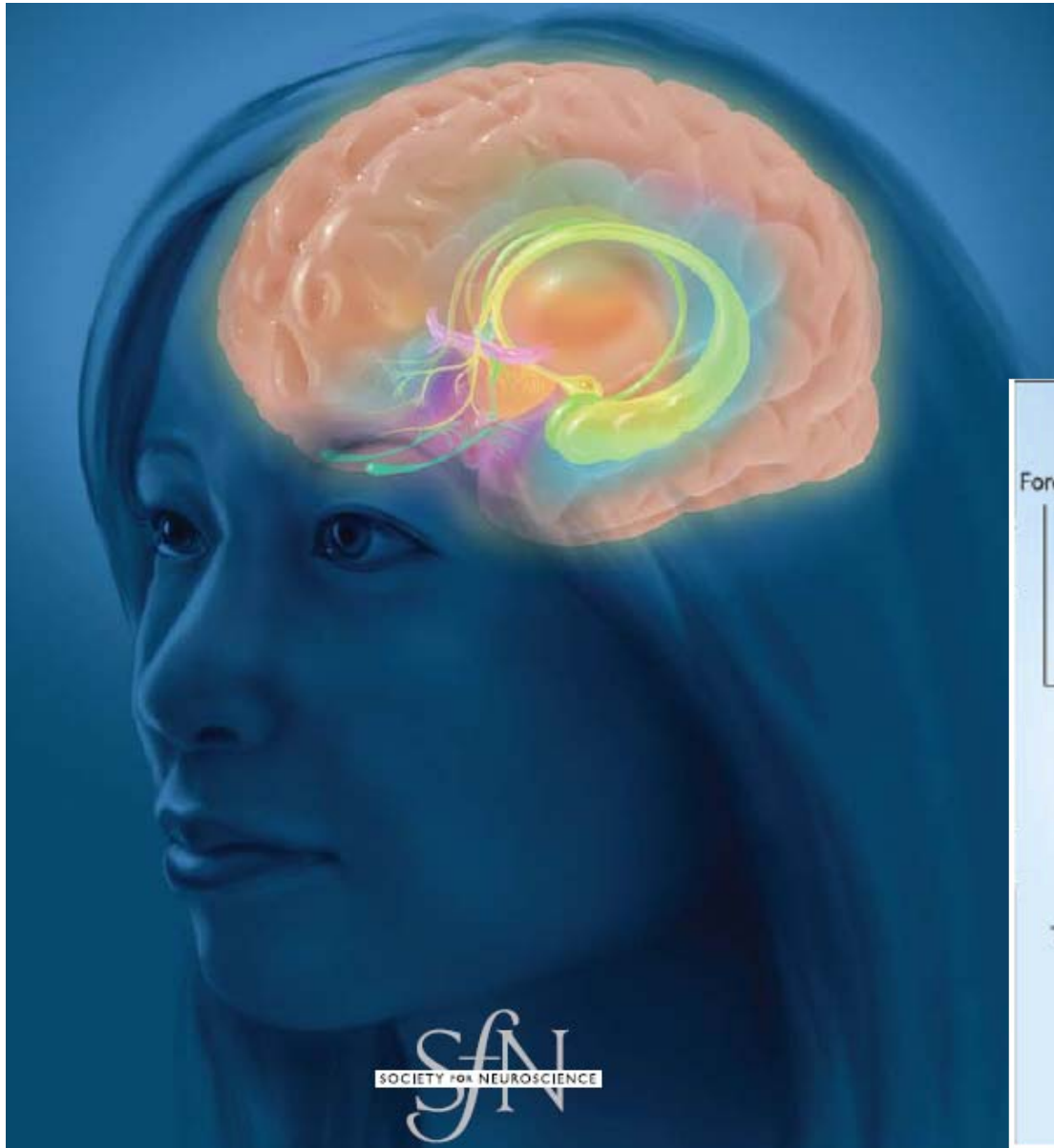
Design Alternatives

➤ Applications

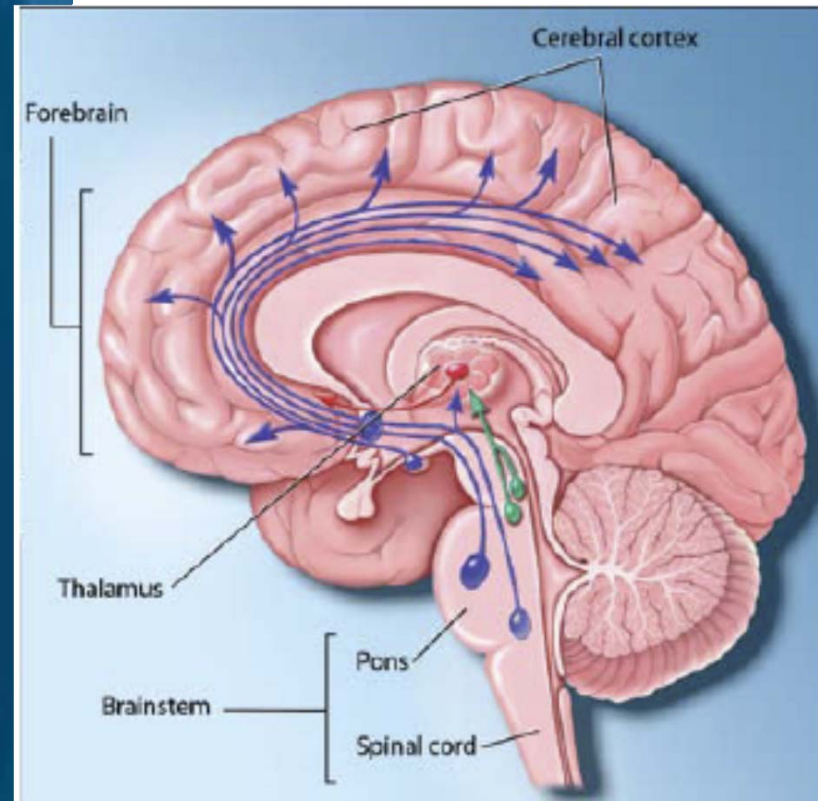
Cognitive Robotics

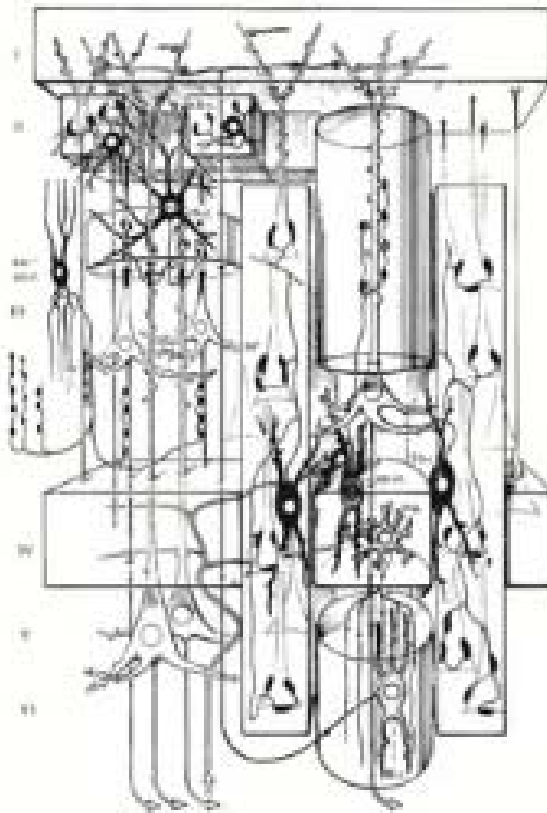
➤ Discussion

Questions

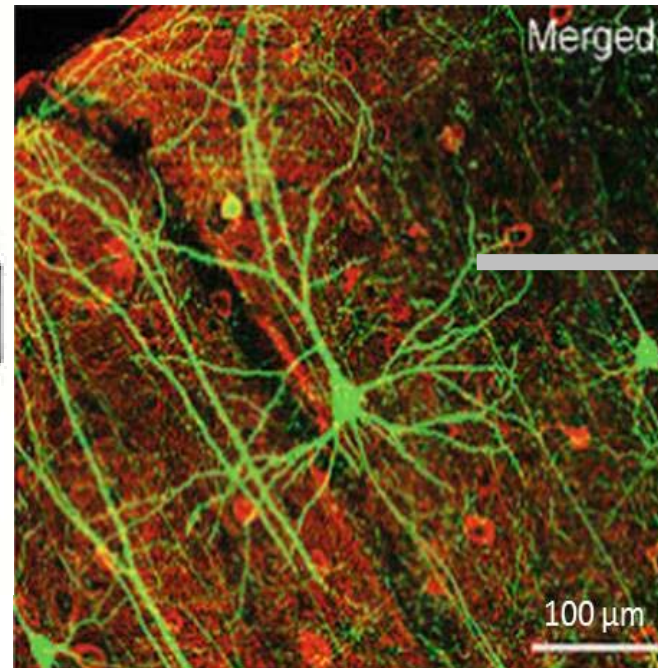


- ~ dm³
- ~ 1,5 kg
- ~ 10¹¹ neurons
- ~ 30 Watt

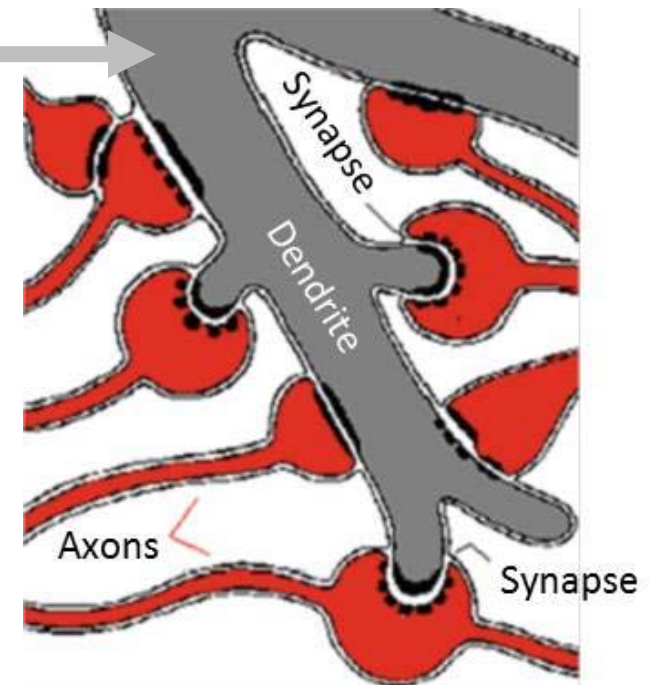




János Szentágothai, *Proc. Roy. Soc. London Ser. B* 201: 219-248)



Synapse
 ~ 0,01 μm^3
 ~ 1ms
 ~ fJ



Cortical column
 ~ 1 mm^3
 ~ 100.000 neurons

Neuron
 < 100.000 μm^3
 ~ 10 Hz (1-100 ms)
 ~ pJ

<http://faculty.washington.edu/chudler/synapse.html>

Architecture

Behaviour

System

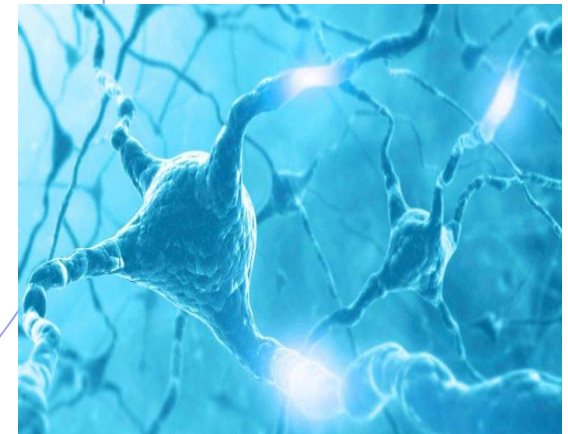
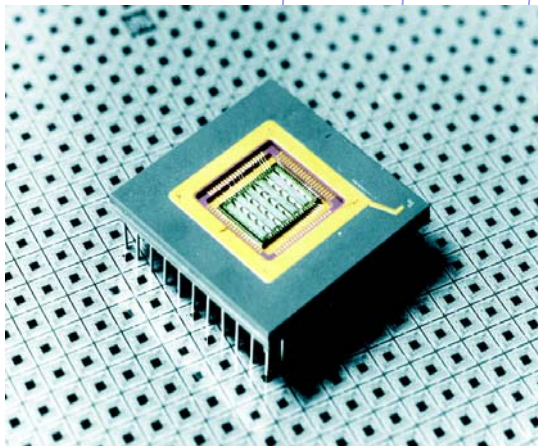
Subsystems

Modules

Logic

Devices

Abstraction Levels



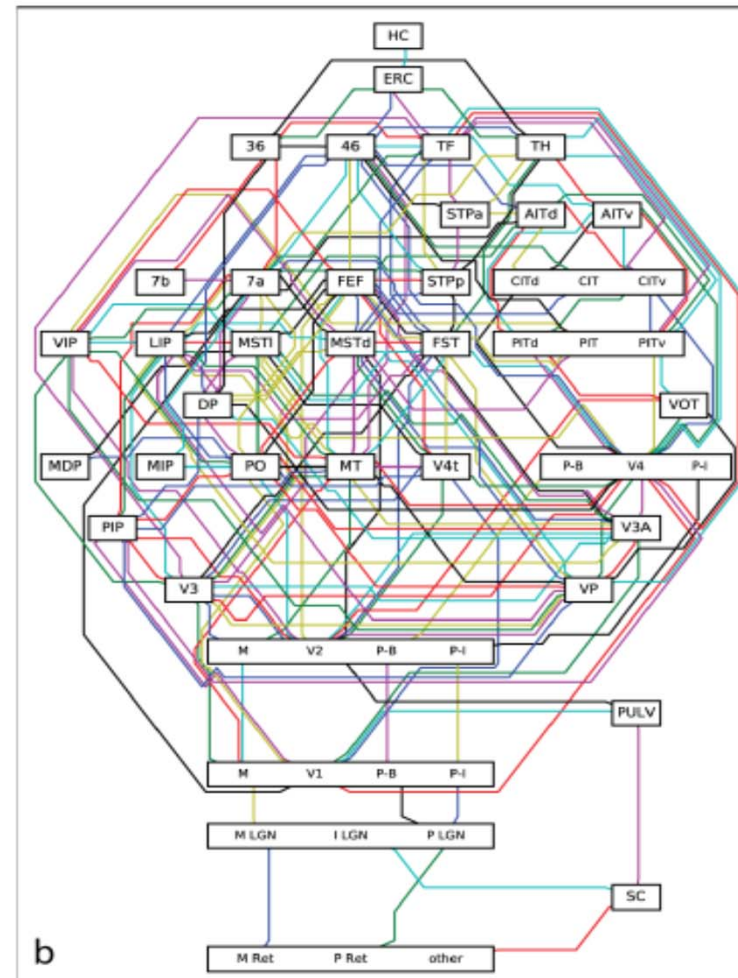
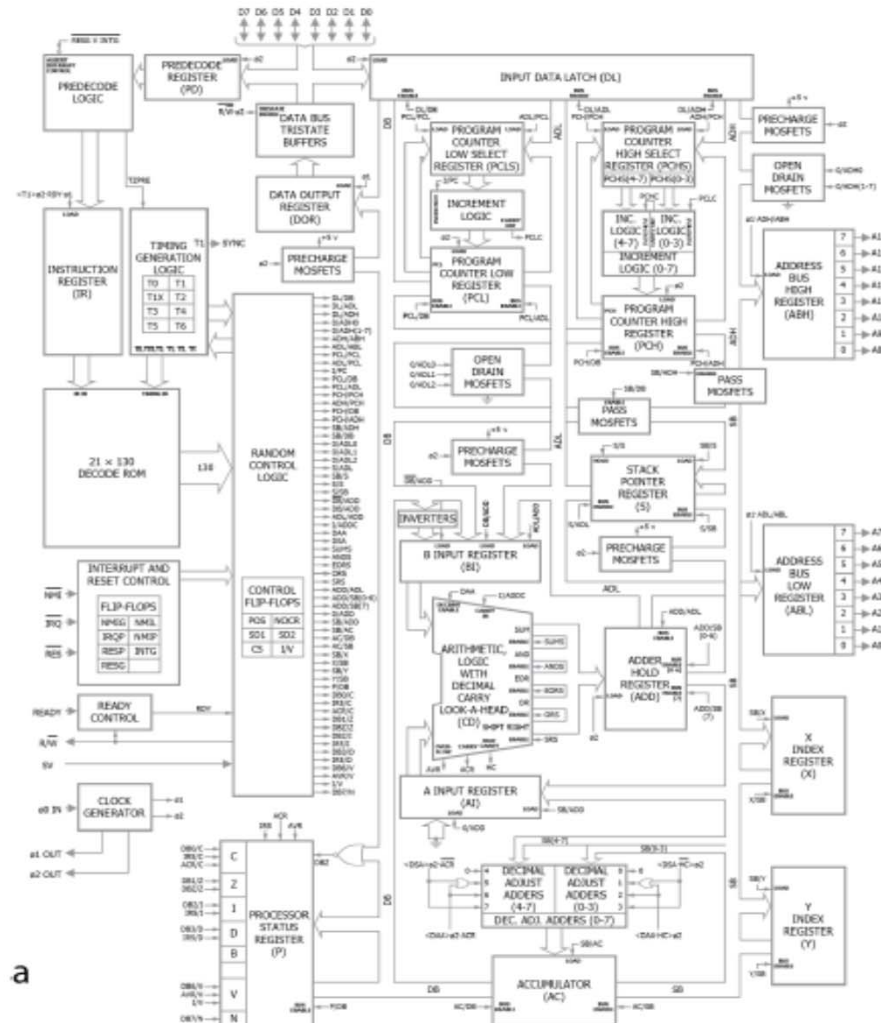
KTSDSIGN Science Photo Library Getty Images

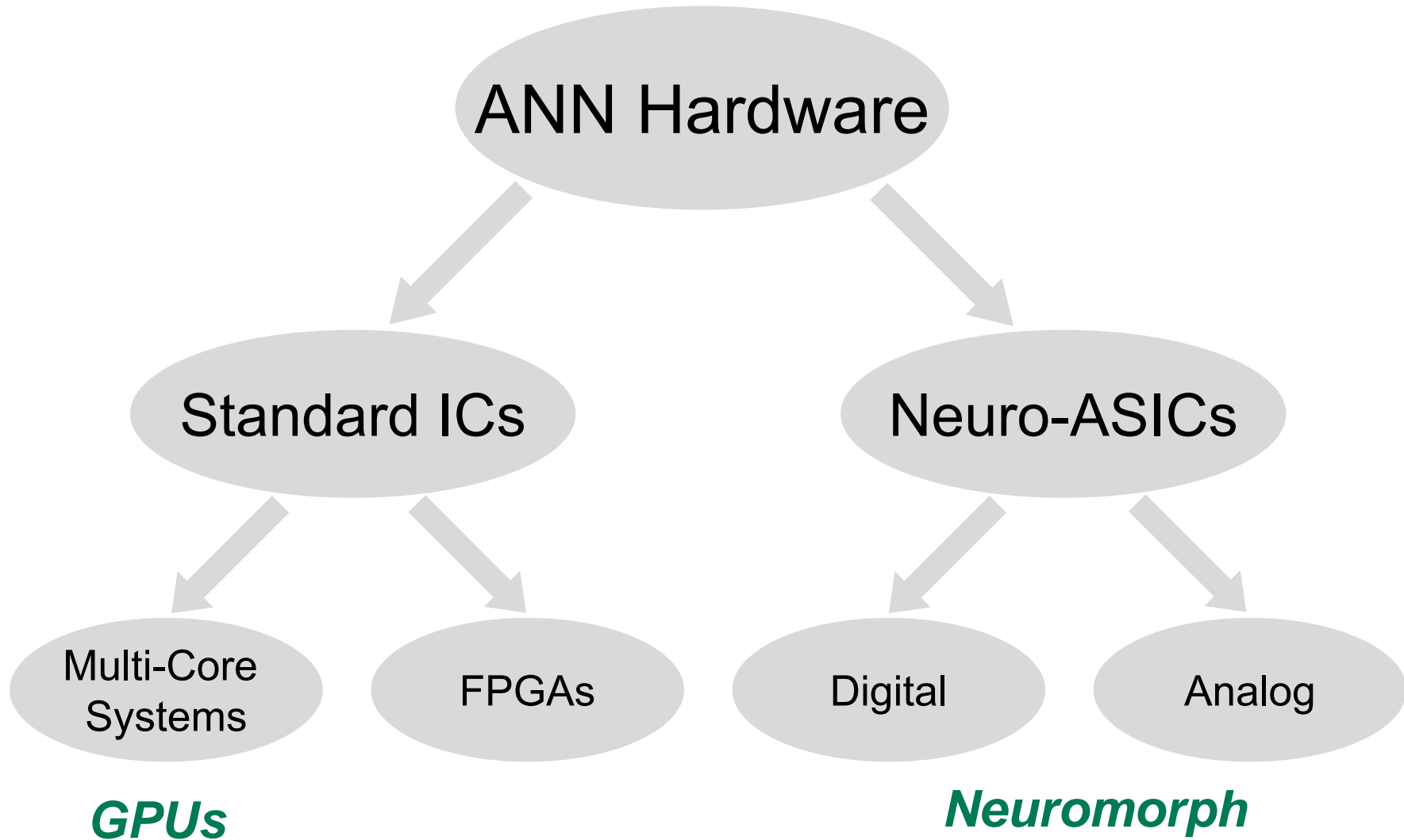
Technology

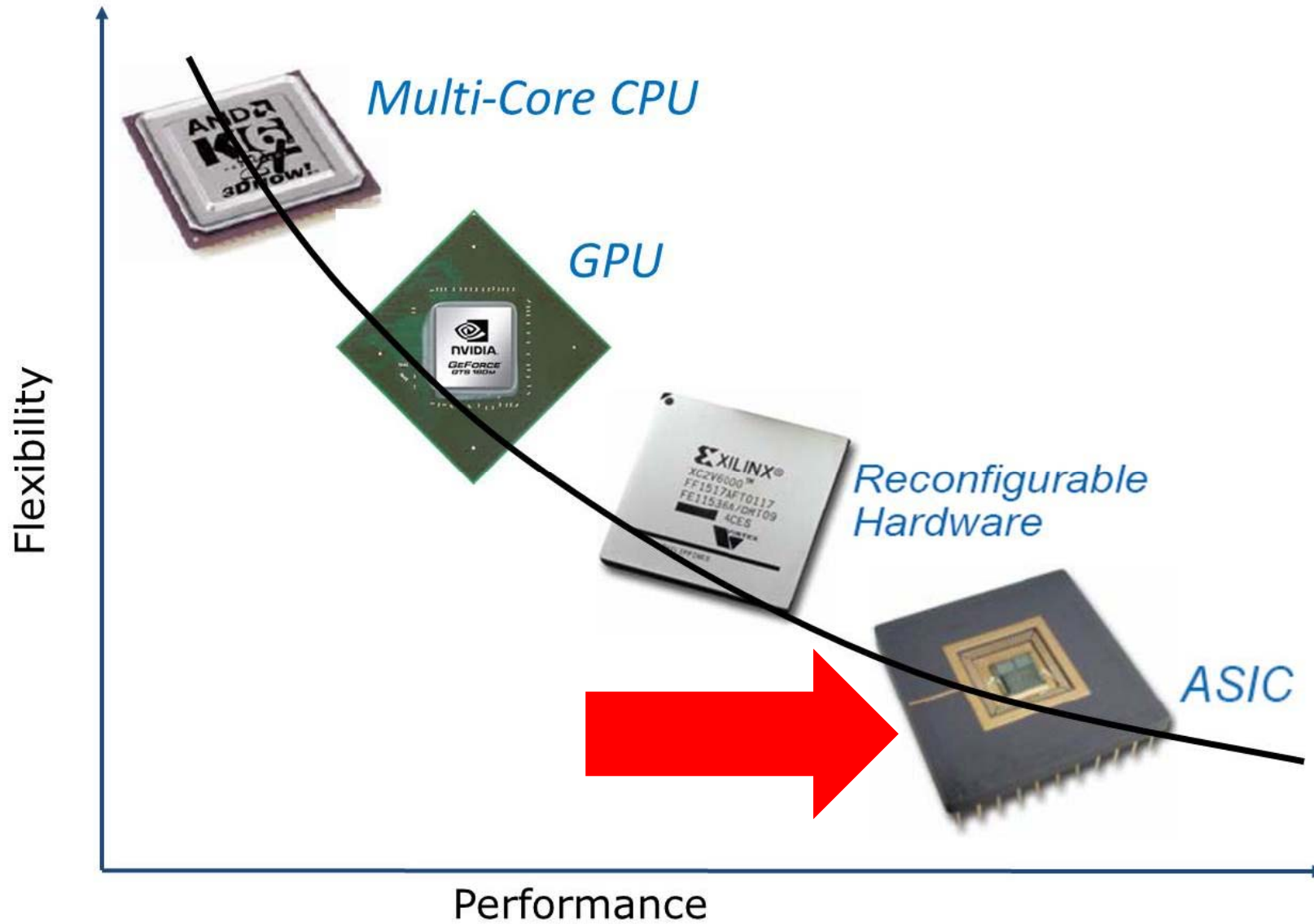
Gajski & Kuhn

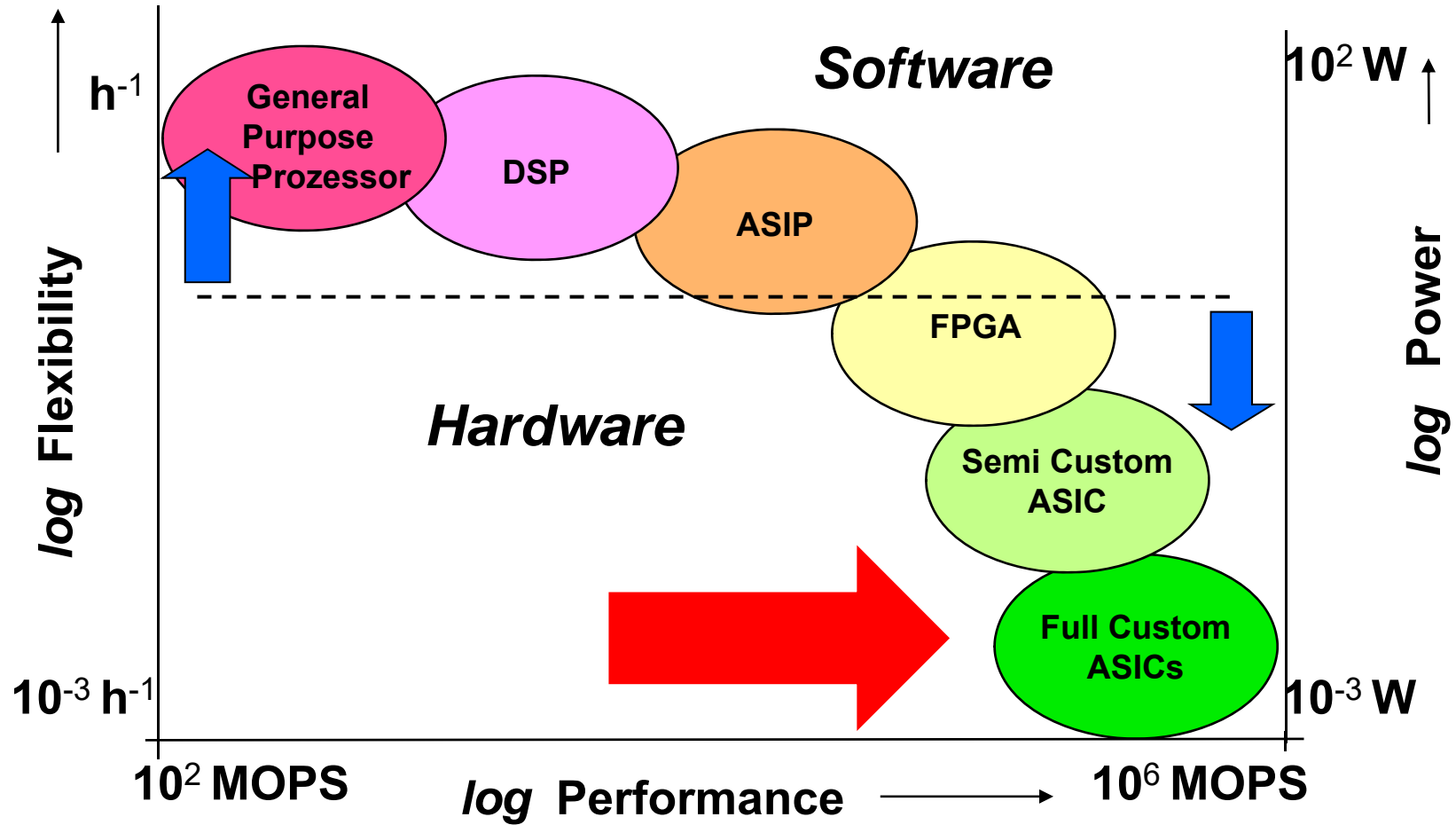
Could a Neuroscientist Understand a Microprocessor?

E. Jonas, K.P. Kording, PLOS May, 2016









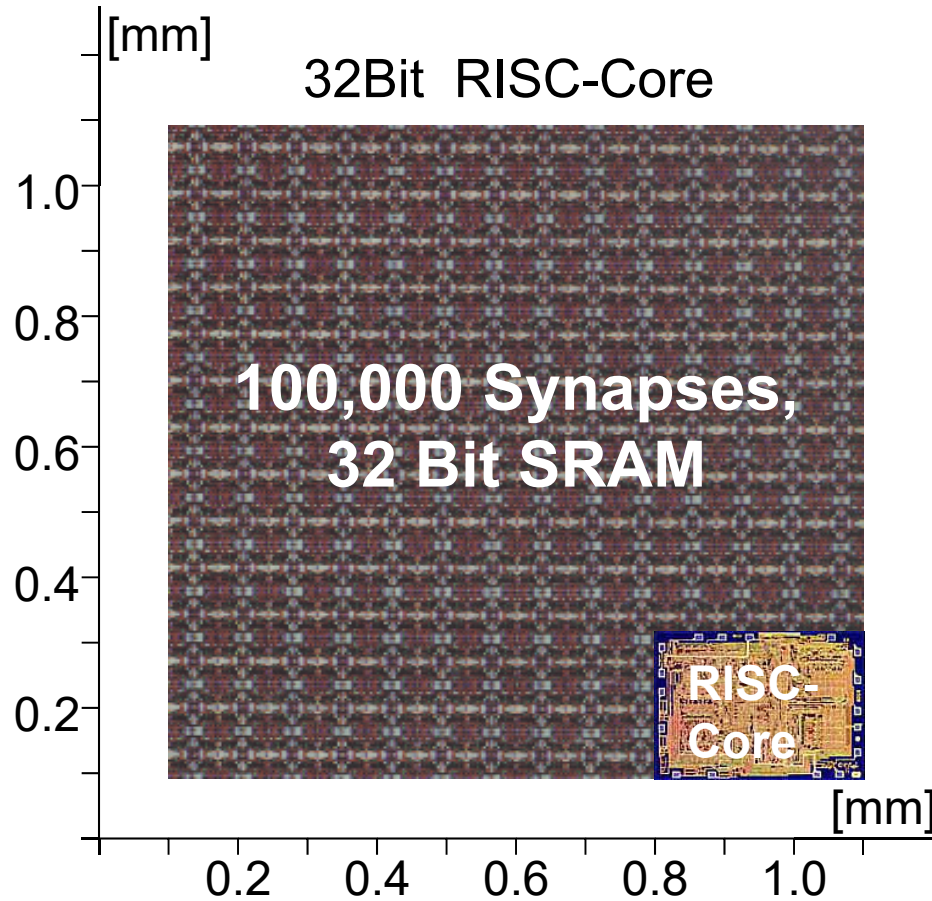
DSP
ASIP

Digital Signal Processor
Application Specific Instruction Set Processor

FPGA
ASIC

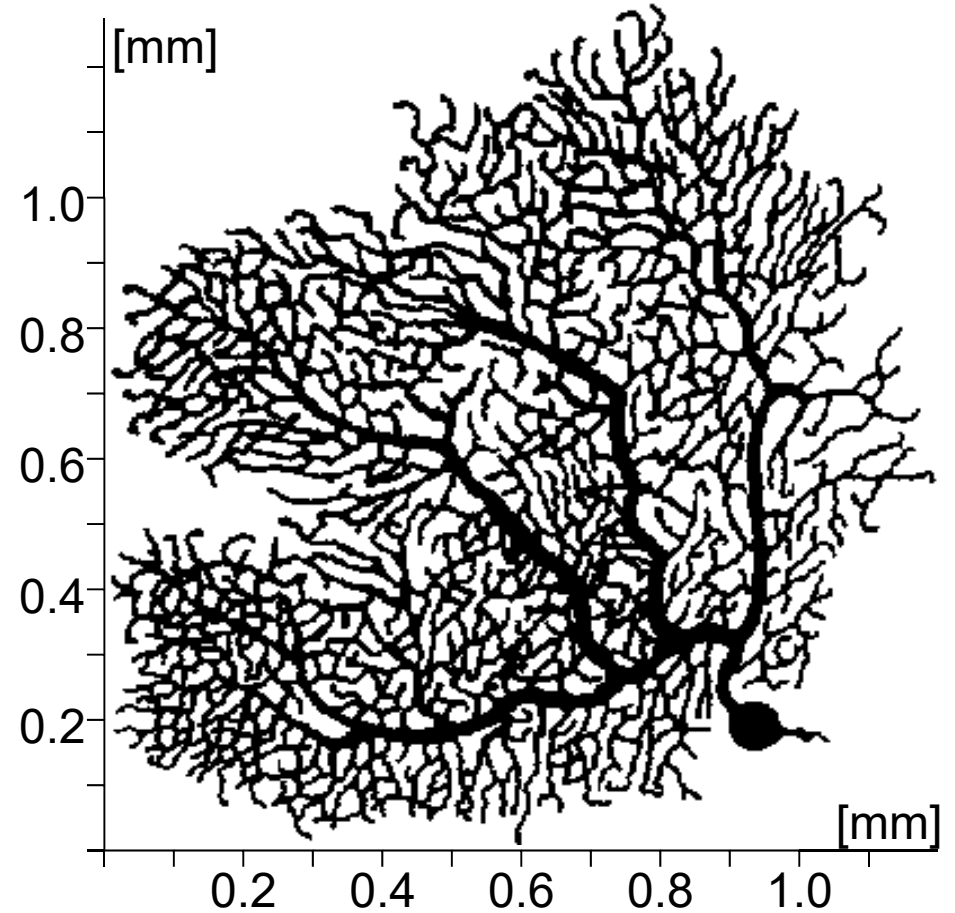
Field Programmable Gate Array
Application Specific Integrated Circuit

Digital-Neuron (μW)



32 nm CMOS

SRAM cell area $\sim 0,1\mu\text{m}^2$



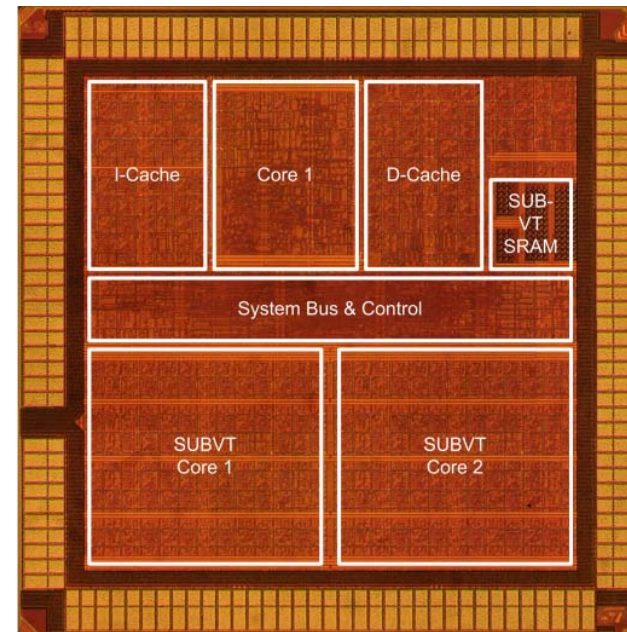
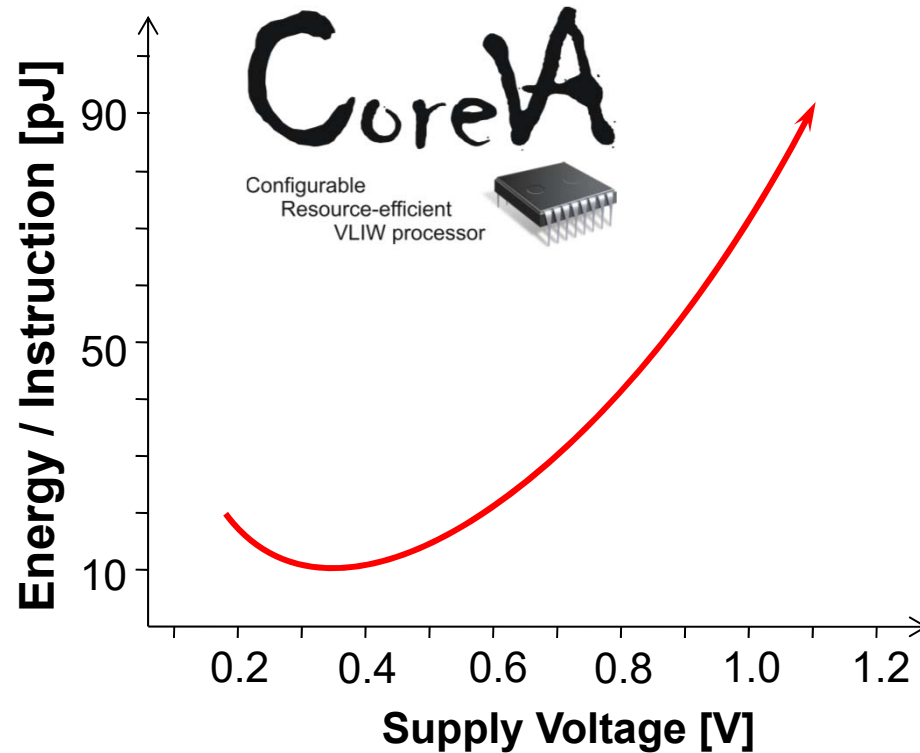
**Purkinje cell of
cerebellar cortex (nW)**

Adaptable supply voltage & operating frequency

Power consumption

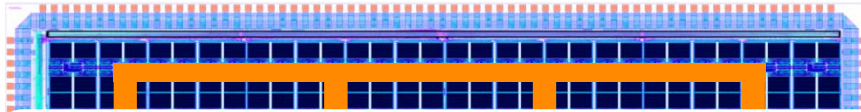
1.35 μ W – 133 kHz – 0.325 V

10,4 mW – 95 MHz – 1.2 V



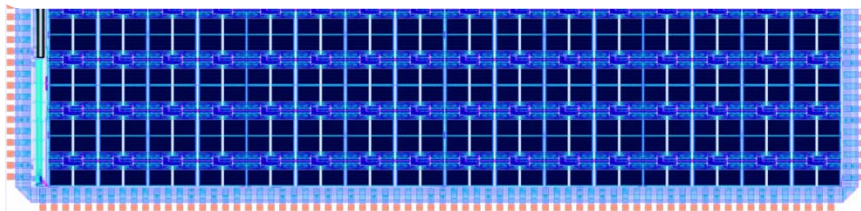
Sub-threshold CMOS
0.36 mm² area, 65 nm

ISSCC2012



Power Challenge:

1000W versus 1W!



1 cm²

Bracaloni Cluster Chip Architecture

- First Chip Realization
 - 2 Cluster / 8 Cores
- 65 nm CMOS Technology
 - 16 Cluster / 64 Cores
- 32 nm CMOS Technology
 - 64 Cluster / 256 Cores
- 16 nm CMOS Technology
 - 256 Cluster / 1024 Cores

It roughly estimates the attainable speedup of a parallel algorithm.

Each algorithm has a parallel and a sequential portion.

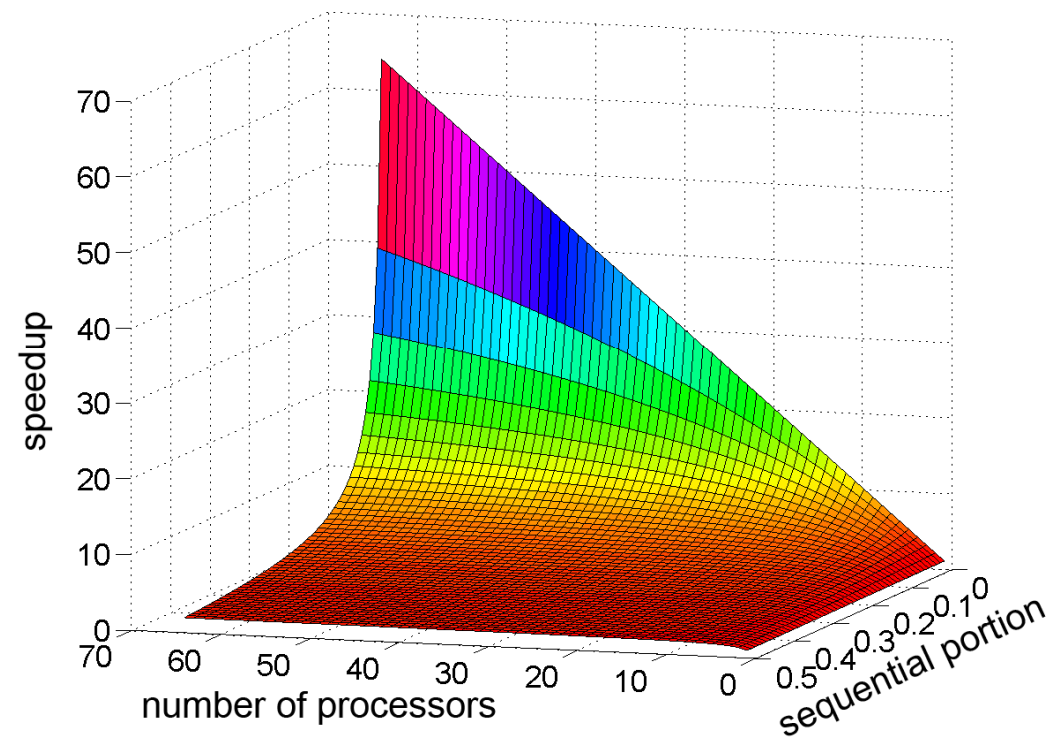
- $\tau_{\text{par}} + \tau_{\text{ser}} = 1$

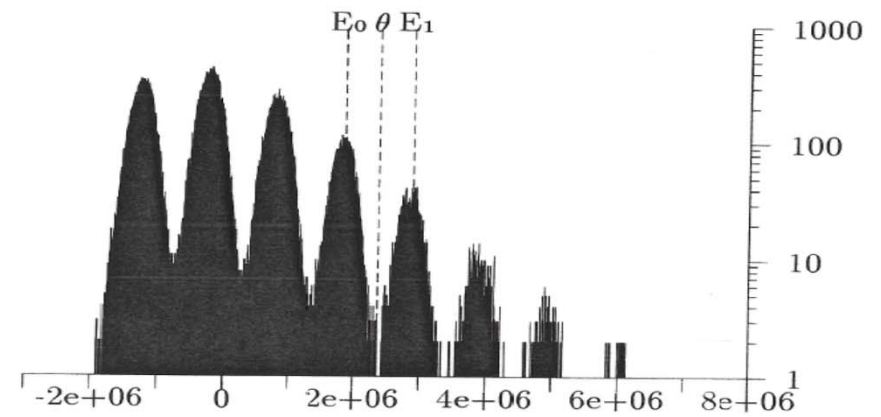
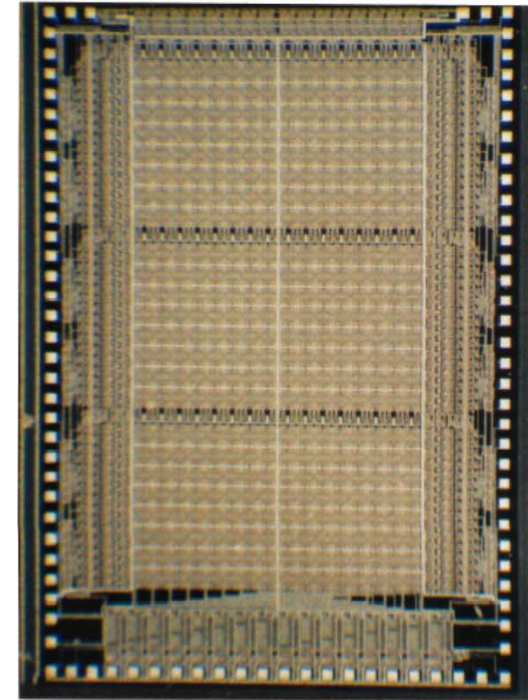
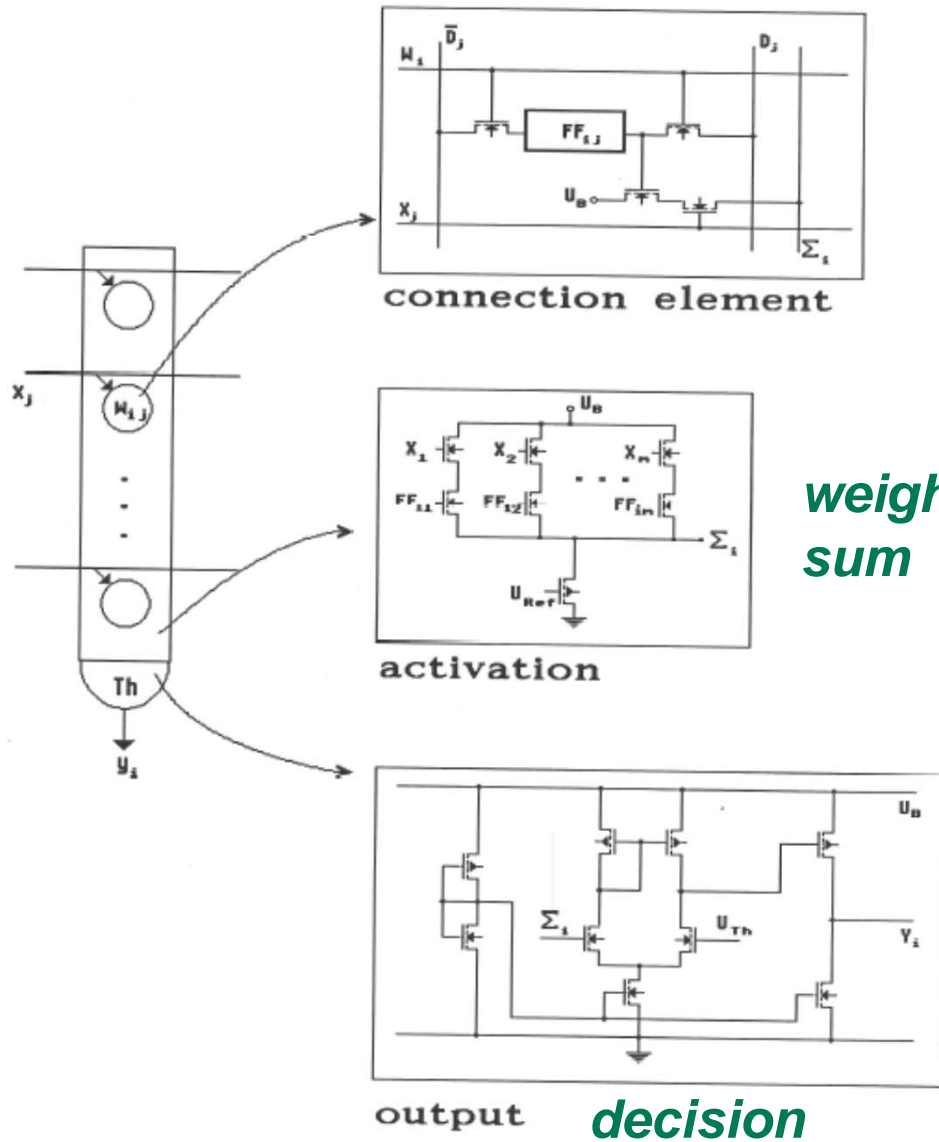
- Runtime with p processors:

$$t(p) = \frac{\tau_{\text{par}}}{p} + \tau_{\text{ser}}$$

- Speedup

$$s(p) = \frac{1}{\frac{\tau_{\text{par}}}{p} + \tau_{\text{ser}}}$$



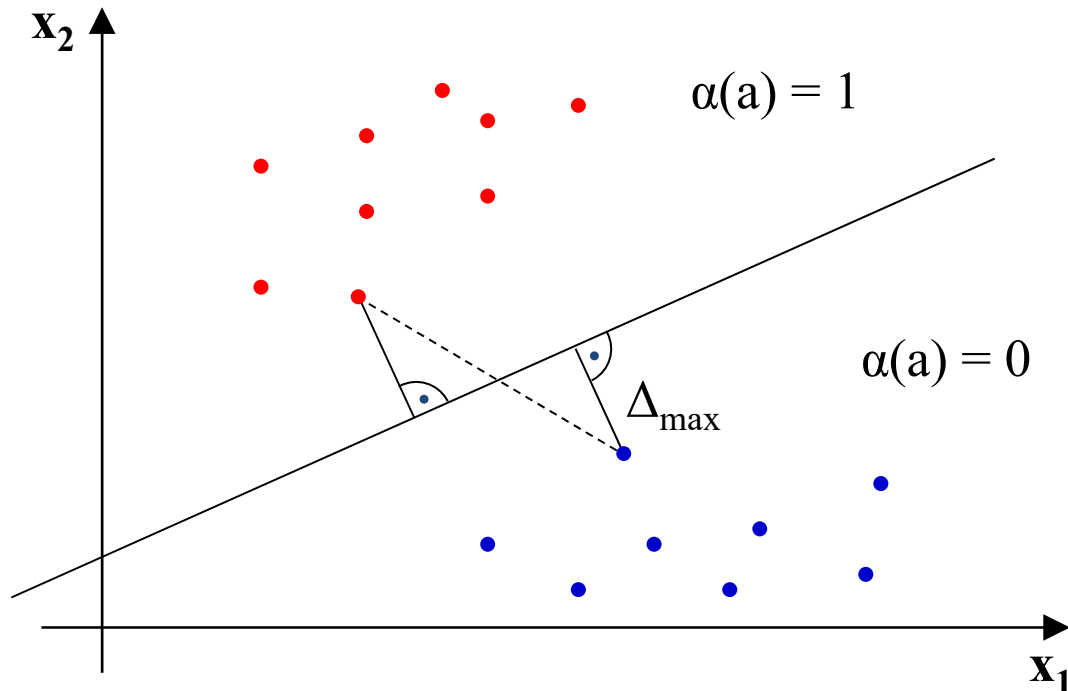


Binary Threshold Neuron:

weights w binary
 Input vector x binary

m = length of x
 l = number of active inputs in x

$$\Delta_{\max} < \min \left\{ \frac{\left| \sum_{i=1}^m w_{ij} \cdot x_i^t - Th \right|}{\sum_{i=1}^m |w_{ij}| \cdot x_i^t + |Th|} \right\}$$



$$\Delta_{\max} < \frac{0.5}{2l - 0.5} \approx \frac{1}{4l}$$



Sparse Coding !

Spiking Neurons

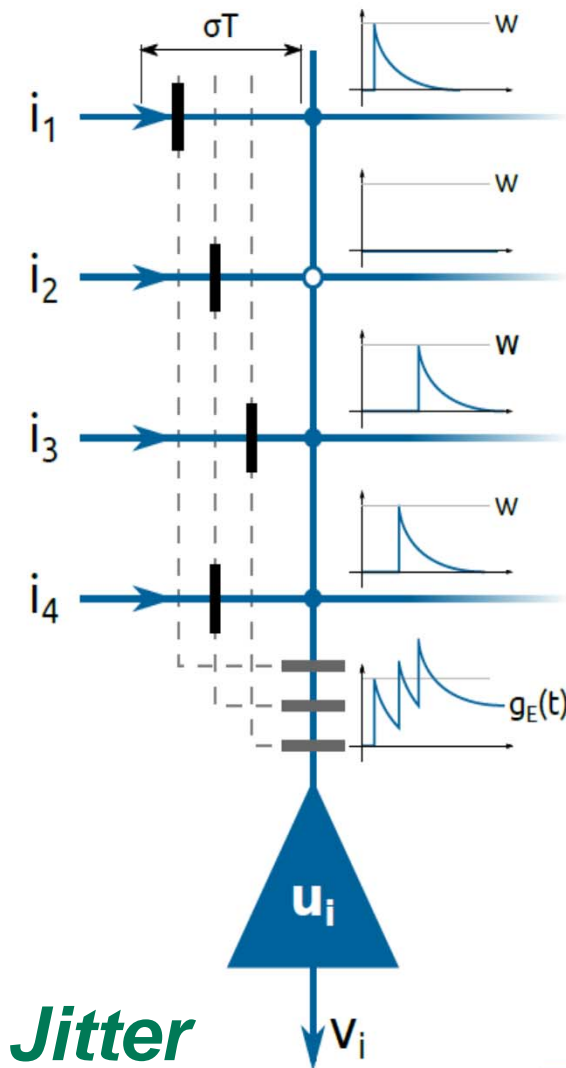
► Given:

- Presynaptic excitatory spikes, jitter σT
- Binary synapse weights:
 $w_0 = 0$ (not connected), $w_1 = w$

► Goal: *(Optimization Problem)*

Find parameters \mathcal{P} with

- N input spikes $\Rightarrow P(\text{spike})$ large
- $N - 1$ input spikes $\Rightarrow P(\text{spike})$ small



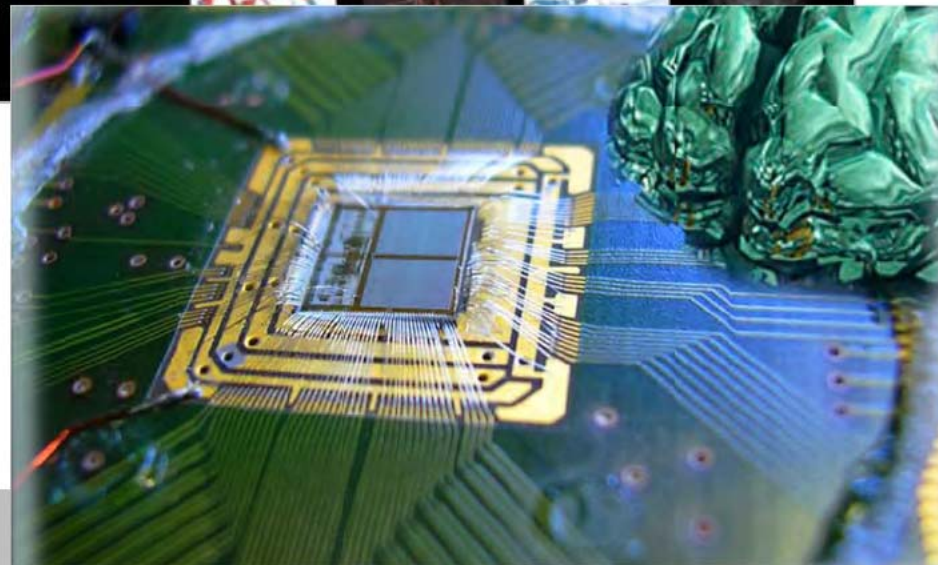
Asynchronous Behaviour \rightarrow Spike Jitter



Two major goals of the Human Brain Project



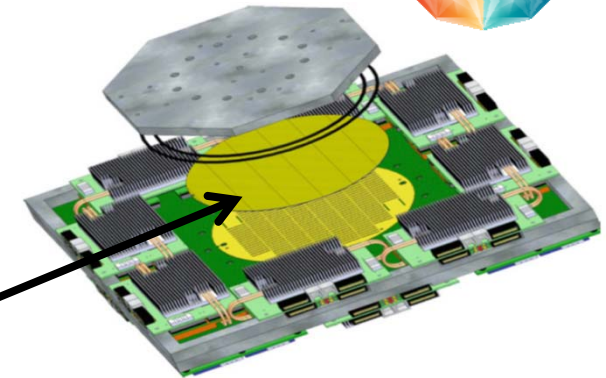
2. Developing powerful simulations, both in conventional computers and in specialized hardware



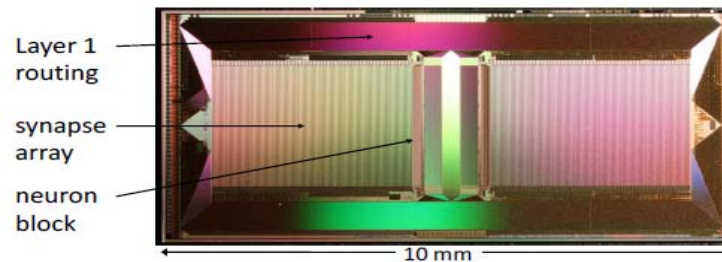


Three (neuromorphic) computing platforms:

- High Performance Computers (HPC)
- Wafer-Scale-System (Heidelberg)

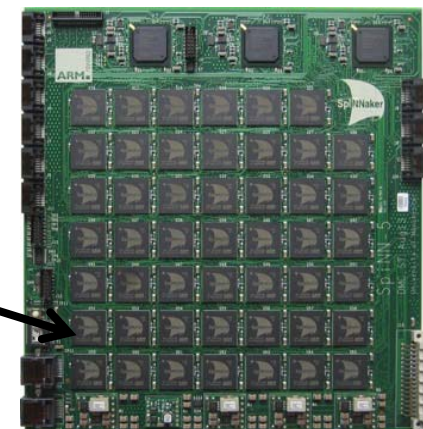
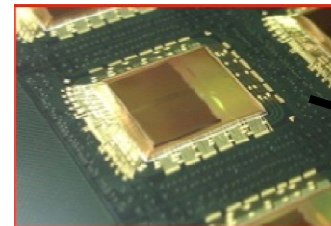


analog
Neurons



- Multi-Core-System (Manchester)

digitally emulated
Neurons





HBP Neuromorphic Computing Concepts



MANY-CORE NUMERICAL MODEL SYSTEM

0.5 – 1 Million ARM processors – address-based, small packet, asynchronous communication – running at real-time

Location : Manchester (UK)

PHYSICAL MODEL SYSTEM

Local analog computing with 4 Million neurons and 50 Million synapses – binary, asynchronous communication – emulation speed is x 10 000 real-time

Location : Heidelberg (Germany)



Neuro-ASICS	Feature Size	Die Size	Neurons	Synapses	Bit/Syn.	ESE
SpiNNaker	130nm	1,02cm ²	1,600	*128x10 ⁶	8	10 ⁻⁸ J
TrueNorth	28nm	4,30cm ²	10 ⁶	256x10 ⁶	1	10 ⁻¹¹ J
HICANN	180nm	0,50cm ²	8-512	114,688	4-8	10 ⁻¹⁰ J
Neurogrid	180nm	1,68cm ²	65,536	*16x10 ⁶	#13	10 ⁻¹⁰ J
numbers per chip				*off-chip	#shared	ESE: Energy/synaptic event

Human Brain: about 10¹¹ Neurons, 10¹⁵ Synapses, 30W average power

- CNAPS, Adaptive Solutions, USA, 1990 (Digital, ASIC)
- Synapse, Siemens, Germany, 1990 (Digital, ASIC)
- ETANN, INTEL, USA, 1990 (EEPROM, Analog, ASIC)
- MoNA, (mixed signal), Porrman/Rückert
-
- SpiNNaker, Univ. of Manchester, U.K., 2007 (Digital, ASIC)
- HICANN, Universität Heidelberg, 2010 (Analog/Digital, ASIC)
- Neurogrid, Stanford University, 2007 (Analog/Digital, ASIC)
- Neurosynaptic Core TrueNorth, IBM, USA, 2010 (Digital, ASIC)
- BlueHive, Cambridge, U.K., 2012 (FPGA)
- Cadence, Tensilica, Google, ARM,
- ...

➤ Introduction

Background

➤ Technology

Femtoelectronics for ANN

➤ Architectures

Design Alternatives

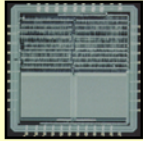
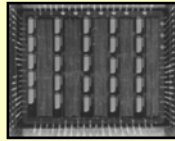
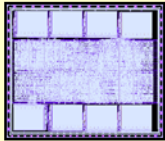


➤ Applications

Cognitive Robotics

➤ Discussion

Questions

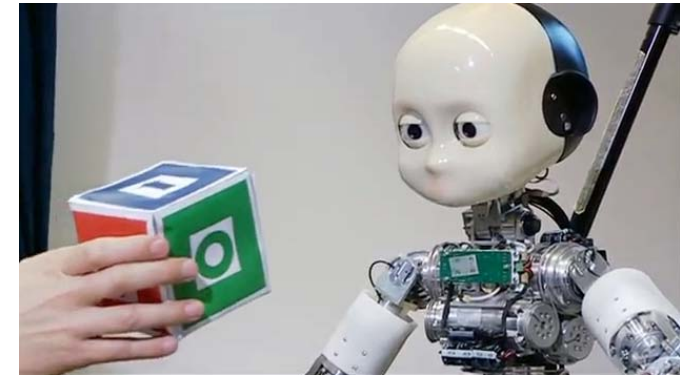
IN → ANN → Out		Task	Evaluation	Model	Realization
discrete	discrete	association	storage efficiency	NAM	
continuous	discrete	classification	error probability	SOM	
continuous	continuous	controlling, approximation	distance measure	LCNN	

- ANN = Artificial Neural Networks
- SOM = Selforganizing Maps

NAM = Neural Associative Memory
 LCNN = Local Cluster Neural Nets



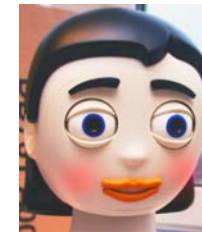
Motion Intelligence



Attentive Systems



Ressource-effiziente
architekturen für
kognitive systeme:
Cognitronics

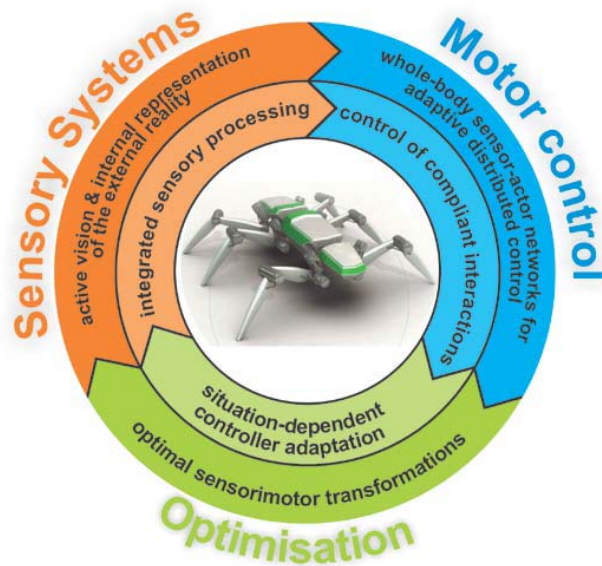


Memory and learning



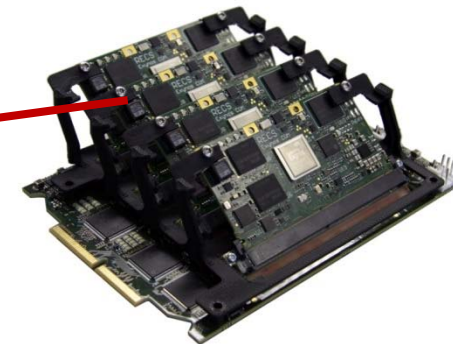
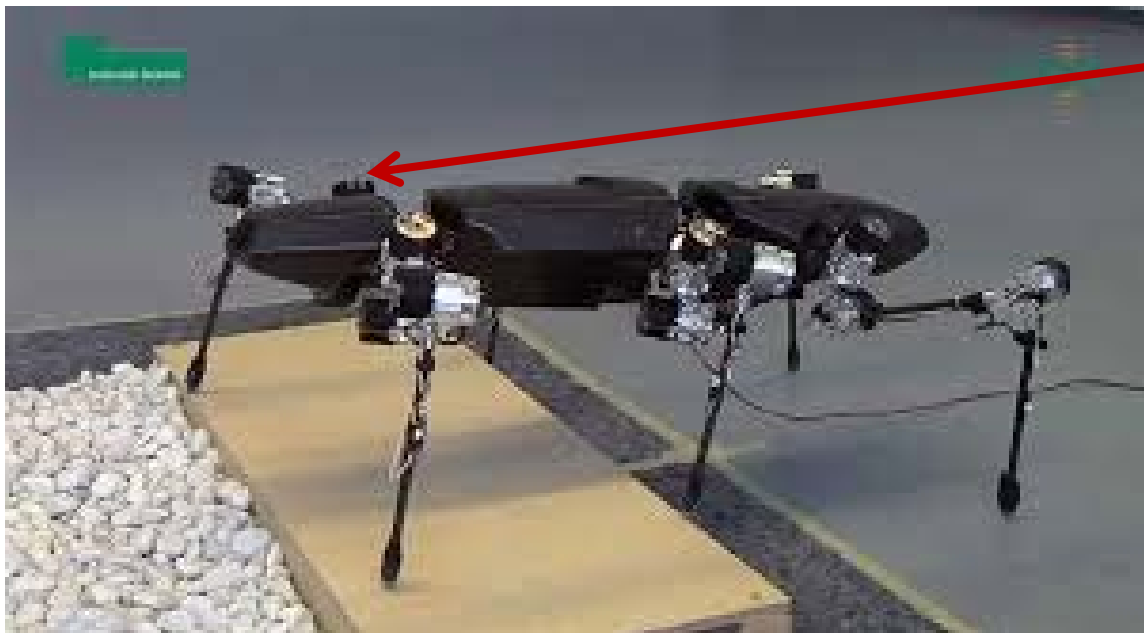
Situated Communication





How to generate complex behaviour in the limits of restricted energy resources?

HECTOR - The six-legged walking robot
see [YouTube](#)

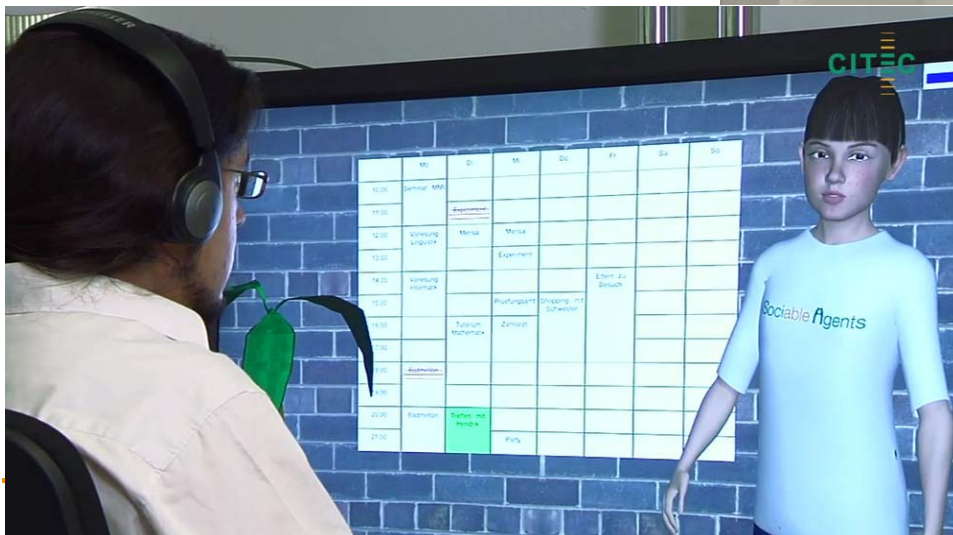


Axel Schneider
 Martin Egelhaaf
 Volker Dürr
 Marc Ernst
 Ulrich Rückert
 Elisabetta Chicca

Long-term experiments in real world scenarios

- 1000 Processors
- TeraByte Memory
- Embedded Sensors

~1dm³
 < 10W



Art Exhibition Hall, Bielefeld

www.cit-ec.de

CITEC Smart Home Apartment



Miele



➤ Introduction

Background

➤ Technology

Femtoelectronics for ANN

➤ Architectures

Design Alternatives

➤ Applications

Cognitive Robotics



➤ Discussion

Questions

- **Sparse codes and acticity (Cell Assemblies)**
reduced power consumption - simple communication
- **Asynchronous information processing**
reduced power consumption - higher robustness
- **Massiv Parallelism and Redundancy**
reduced power consumption - higher robustness
- **Continuous Self-organization**
fault-tolerance - optimal use of resources
- **Stability in large, dynamic networks**
local rules for stable and robust global behaviour

**1 Billion Cores
(Peta Byte)**



System in a Rack



**1 Mill. Cores
(Tera Byte)**



1dm³
100W

System in a Box



**1.000 Cores
(Giga Byte)**



1cm³ ; 1W

System in a Dice





System in a Rack



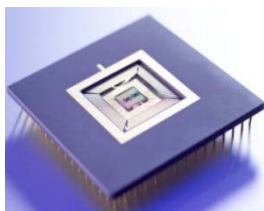
System in a Box



1dm³
100W



System in a Dice



1cm³ ; 1W



At least in the next decade neural network implementations are dominated by nanoelectronics.

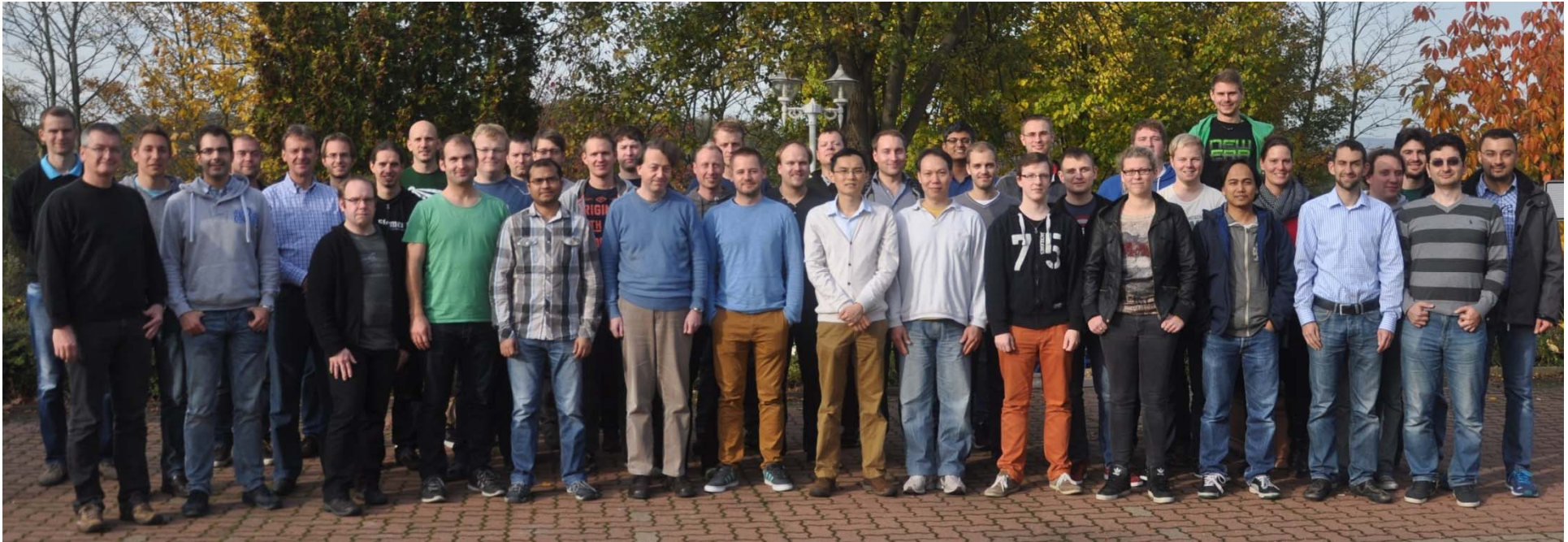
Until fundamental concepts become better understood, the main advantages of silicon should not blind us to alternative neural network design.

E. R. Caianiello



Thank you for your attention!

Ulrich Rückert



The Cognitronics and Sensor Systems Research Group