# **The FACETS Project**

# VLSI Implementations of Very Large Scale Neuromorphic Circuits

#### **Johannes Schemmel**

Electronic Vision(s) Group Kirchhoff Institute for Physics Ruprecht-Karls-Universität, Germany

The Heidelberg FACETS team (in no particular order):

Karlheinz Meier, Sebastian Millner, Dan Husmann, Andreas Grübl, Johannes Schemmel, Maurice Güttler, Holger Zoglauer, Matthias Hock, Simon Friedmann, Stefan Philipp, Daniel Brüderle, Mihai Petrovici, Moritz Schilling, Johannes Bill, Bernhard Kaplan, Erik Müller, Andreas Hartel

## Why Build Artificial Neural Systems?



- we want to understand the biological computing paradigm
  - limited experimental access
  - artificial neural systems act as a model for biology
- new applications → robotics, ambient intelligence, user interfaces, etc

## **Basics of Neural Communication**



© brainmaps.org

- neurons integrate over space and time
- temporal correlation is important
- kind of mixed-signal system:
   action potential ↔ membrane voltage
- fault tolerant
- low power consumption →
   100 Billion neurons: 20 Watts



## **Modeling approaches for Artificial Neural Systems**

#### starting point: mathematical description

methods:

- analytical treatment proof of general properties and limits
- numerical solution (high performance computing) flexibility, parallel objects not obvious
- physical model (neuromorphic hardware) artificial nervous system artificial parallel object = biological objects
- biological model

"custom-made biological nervous system"

## **Limits of Numerical Models**



- loss of fault tolerance inherent to neural systems
- power consumption of the simulation layer



biologically inspired architectures preserve the fault tolerance and low power consumption of neural systems at the device level  $\rightarrow$  physical model

#### International Technology Roadmap for Semiconductors (ITRS)



# **The FACETS Project**



#### FACETS

Fast Analog Computing with Emergent Transient States

## FACETS - From Neurobiology to new Computing Architectures



U Bordeaux, CNRS (Gif-sur-Yvette and Marseille), U Debrecen, TU Dresden, U Freiburg, SCCH Hagenberg, TU Graz, U Heidelberg, EPFL Lausanne, U London, U Plymouth, INRIA Sophia-Antipolis, KTH Stockholm

An Integrated Project in the 6th Framework Programme Information Society Technology - Future Emergent Technologies FP6-2004-IST-FETPI



Ruprecht-Karls-Universität Heidelberg

#### FACETS : Basic Idea, methodological approach and goals

**Neurobiology** : Structural and Functional Investigation of the Neocortical Microcircuit and the Circuit Elements in-vivo and in-vitro



Ruprecht-Karls-Universität Heidelberg

Johannes Schemmel

## **Essentials for Neuromorphic Hardware Systems**

#### **Biological Neural Computation**

- Connectivity 10<sup>11</sup> Neurons, 10<sup>15</sup> Synapses in Neocortex 10.000 Synapses per Neuron on average
- Diversity Categories and Parameters of Neurons
- Plasticity Long Term, Short Term, Local, Global
- Timing Time constants, delays, correlations

#### Neuromorphic Hardware Systems

- Connectivity Efficient data protocols, 2D-3D connection technology
- Diversity Configurability (distributed memory)
- Plasticity Local and global dynamic and static memory
  - Control time constants, delays and time correlations
- SCALABILITY Learn from small systems Approach large scales Bandwidth, delays, power, cost, fault tolerance

Timing

Гb

## **Continuous Time Integrating Membrane Model**

Consider a simple physical model for the neuron's cell membrane potential *V*:

	ΔV [V]	$g_{\text{leak}}[S]$	C <sub>m</sub> [F]	(gV)/C [V/s]	Inherent speed gap:
Biology(*)	10-2	10-8	10-10	100	10 <sup>6</sup> Volt/second
VLSI	10-1	10-6	10-13	106	→ <u>accelerated</u> <u>neuron model</u>

(\*) from Brette/Gerstner, J. Neurophysiology, 2005

# **Properties of the Facets Neuron Model**

#### Accelerated continuous time:

- acceleration factor 10<sup>3</sup> to 10<sup>5</sup>
- high speed operation for statistics and parameter searches
- minimum area leads to high synapse numbers
- No deep sub-threshold operation in principle all analog parameters could be well controlled → automatic transfer of experiments (PyNN)

## FACETS Stage 1 : Conductance-based Network Model



#### Synapses:

- $p_{k,l}(t)$  exponential onset and decay (spike shape)  $g_{k,l}$  0 to  $g_{max}$  with 4 bit (8 bit) resolution
- effective membrane time-constant  $c_m/g_{total}$  is time-dependent

#### Stage 2: FACETS Adaptive-Exponential Integrate-and-Fire Neuron Model



R. Naud et al.: Versatility and Relevance of the adaptive-exponential integrate-and fire-model, Biol Cybern (2008) 99:335–347

Ruprecht-Karls-Universität Heidelberg

#### **CMOS VLSI Implementation of the AdExp Integrate-and-Fire Neuron**





## Spike Firing Modes of the AdExp VLSI Neuron

**Transient Response** 1.2 spike frequency adaptation 1.15 1.1 ۔ ج > .95 .9 .85 1.2 phasic spiking 1.15 1.1 2 2 2 2 2 .95 .9 .85 1:2<sup>-1</sup> tonic spiking L 1.15 1.1 2<sup>1.05,</sup> 2, > .95 .9 .85 600 stimulus 500 400 € \_\_\_\_\_300 200 100 0 80.0 90.0 100 110 120 130 140 time (us)

Ruprecht-Karls-Universität Heidelberg

Johannes Schemmel

#### Kirchhoff Institute for Physics

## **Burst Firing Modes of the AdExp VLSI Neuron**



Ruprecht-Karls-Universität Heidelberg

# **Stage 1 : Chip Specifications**

- technology: UMC 180 nm, 6 metal layers, 1 polysilicon layer
- chip size: 5 x 5 mm<sup>2</sup>
- 384 neurons, 100k synapses
- scale factor up to 100k : 10 ns chip-time equals 1 ms real-time
- fast analog outputs (about 400 MHz bandwidth) to monitor selected membrane potentials
- internal storage for model parameters (about 4k values)



A New VLSI Model of Neural Microcircuits Including Spike Time Dependent Plasticity, Johannes Schemmel, Karlheinz Meier, Eilif Muller, Proceedings of the 2004 International Joint Conference on Neural Networks (IJCNN'04), IEEE Press, pp. 1711-1716, 2004

# **Stage 1 : Circuit Features**

- fully analog network core
- continuous time network operation
- short-term synaptic depression and facilitation: analog on-chip
- Spike Time Dependent Plasticity measurement in each synapse, weight update performed digitally
- programmable model parameters (individually or group-wise):
  - reversal potentials: excitatory, inhibitory and leakage  $(E_x, E_i, E_l)$
  - threshold voltage level V<sub>th</sub> and comparator speed
  - reset potential ( $V_{\text{reset}}$ ) and leakage conductance ( $g_{\text{leak}}$ )
  - synapse parameters: rise time, fall time, maximum conductance ( $t_{rise}$ ,  $t_{fall}$ ,  $g_{k,l \max}$ )

A New VLSI Model of Neural Microcircuits Including Spike Time Dependent Plasticity, Johannes Schemmel, Karlheinz Meier, Eilif Muller, Proceedings of the 2004 International Joint Conference on Neural Networks (IJCNN'04), IEEE Press, pp. 1711-1716, 2004

Ruprecht-Karls-Universität Heidelberg

#### **Response to Poisson-Distributed Input Spike Trains: Hardware vs. Simulation**



## Stage 1 Setup



Ruprecht-Karls-Universität Heidelberg

#### **Example: Study of self-stabilizing network architectures**

- Based on short-term synaptic plasticity, to counterbalance transistor level hardware variations
- Various parameter sets, multiple repetitions with new random connections each
- Once with static, once with dynamic synapses
- Result: Stable average firing rates, increased run-to-run reliability



#### FACETS mixed-signal VLSI System Stage 1 (Chip based)



Ruprecht-Karls-Universität Heidelberg

Johannes Schemmel

Kirchhoff Institute for Physics

## Solution: The FACETS Wafer-Scale System



Ruprecht-Karls-Universität Heidelberg

#### **Connection Requirements**



- biological connection density of 1000-10000 inputs per neuron
- only two dimensions for wiring (the brain uses three)
- acceleration factor of 10<sup>4</sup>, biological event rate 20 Hz
- $\rightarrow$  resulting signal density: approx. 10<sup>10</sup> events/(mm\*s)

## **Routing Properties**

spatial communication bandwidth 10<sup>10</sup> events/(mm\*s)
 this seems possible, but
 events must be routed from synapse A to synapse B

routing properties:

- constant propagation delay
- routing topology changes only slowly with time (long term plasticity)
- programmable topology for different experiments and neural circuits

solution: one connection per synapse → connection-based routing

resulting connection density: 10<sup>4</sup> connections/mm

not possible with traditional bonding techniques

 $\rightarrow$  wafer scale integration allows a connection pitch close to  $5\mu m$ 

remaining gap: use time multiplexing

#### **Overview of the Wafer-Scale System**



## **Wafer-PCB Alignment and Assembly Facility**



- Large (approx.10.000) number of vertical elastomeric contacts demand precise alignment between wafer and PCB with the peripheral electronics / power delivery. Pitch 200 μm.
- Adjustment accuracy better than 50 μm over the 8 inch wafer

#### **Post Processing Results**

- The measured yield for 8 μm pitch test structures was 100%
- The pitch needed for the final system is 9.5 μm
- With larger pad separation even smaller structures could be realized



Microscopic views of an 8  $\mu m$  pitch post processing structure

Ruprecht-Karls-Universität Heidelberg

Johannes Schemmel

problem: communication bottleneck between wafers consider 10Hz biological firing rate at a an acceleration factor of  $10^4$ assume we implement about  $10^5$  synapses on a 5x10 mm<sup>2</sup> die  $\rightarrow 10^{10}$  action-potentials/s

imagine an AER protocol transmitting the target synapse number (14 bit)
→ about 100 Gbit/s (very optimistic considering the collision probability)

idea: a neighboring chip needs about the same pre-synaptic signals

- → imagine a fast connection between two chips: each chip would need only half the bandwidth to the rest of the world
- $\rightarrow$  a wafer: 350 chips, external bandwidth reduced to GBits/s

# Digital Network Chip (DNC, TU Dresden)



## **Structure of a FACETS Wafer**



## **Low-Power Neuron to Neuron Communication**



## **Power Consumption of the Serial Repeater Circuits**



## High Input Count Analog Neural Network (HICANN)



## High Input Count Analog Neural Network (HICANN)



combining multiple membrane circuits with 256 synapses each allows neurons with up to 16k pre-synaptic inputs

→ will be extended to real multi-compartment neurons with back-propagating actionpotentials by current PhD

### Microphotograph of the HICANN Die



#### 10 mm

Ruprecht-Karls-Universität Heidelberg

#### **Adding Interaction to the System**



latency of host

- communication can be low:
  - O(μs) for GBit-Ethernet
  - <100ns for specialized High-Speed serial links

#### possibilities:

- sensor-motor interactions: movement of a simulated retina
- robot in virtual environment simulated on PC cluster

#### $\rightarrow$ required latency 10 to 100ms: 1 to 10 $\mu s$

## **Software : From Networks to Experiments**

#### **PyNN script**



#### Configuration/Evaluation (comparing connection matrix)





# Summary

FACETS has developed solutions to enable large scale analog hardware as an alternative to numerical simulations:

- programmability of topology and model parameters
- flexible, biologically realistic neuron model
- low-power continuous-time communication
- scalable packet-based inter-wafer and host communication
   → possibility of interactive simulations (sensor-motor loop)
- a software framework for the translation of experiments from biology to hardware