

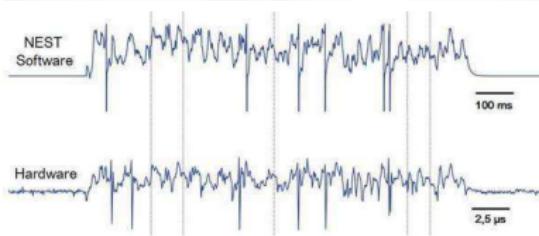
# The PyNN module for the FACETS waferscale neuromorphic hardware system `pyNN.hardware.stage2 (& stage1)`

Eric Müller

<http://www.kip.uni-heidelberg.de/vision/>  
Kirchhoff Institute for Physics,  
University of Heidelberg

2009-10-09

# FACETS Hardware Stage 1



- 384 gLIF neurons,  $\sim 98\text{ k}$  synapses
- highly accelerated operation (speed-up  $10^4 - 10^5$ )
- long- and short-term plasticity
- mapping biological  $\Leftrightarrow$  hardware value domain
- remote access, configurable neuron and synapse parameters
- measurable: spike output, membrane potential, weights

Schemmel et al.: Modeling Synaptic Plasticity within Networks of Highly Accelerated I&F Neurons, ISCAS 2007

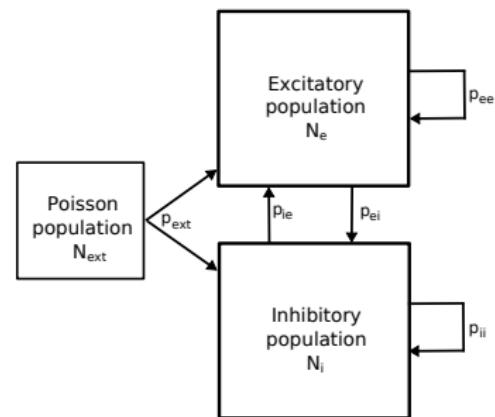
# Software status

- `pyNN.hardware.stage1:`
  - procedural PyNN interface
  - object-oriented PyNN interface
  - simple check of hardware constraints
- `up & running`
- `experiments ...`

## Recurrent network dynamics

# Recurrent network dynamics

- Daniel Brüderle<sup>a</sup> & Jens Kremkow<sup>b</sup>
- compare the Stage 1 hardware system to a simulator
- systematical parameter sweep to generate different states of network activity
- e.g. this recurrent network architecture

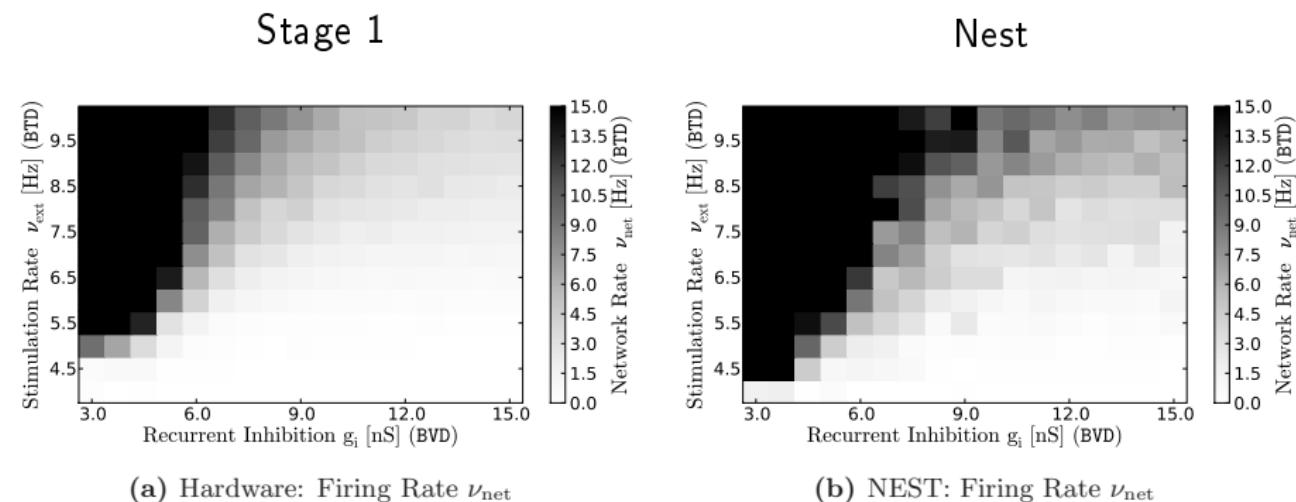


<sup>a</sup>University of Heidelberg

<sup>b</sup>CNRS, Marseille, France

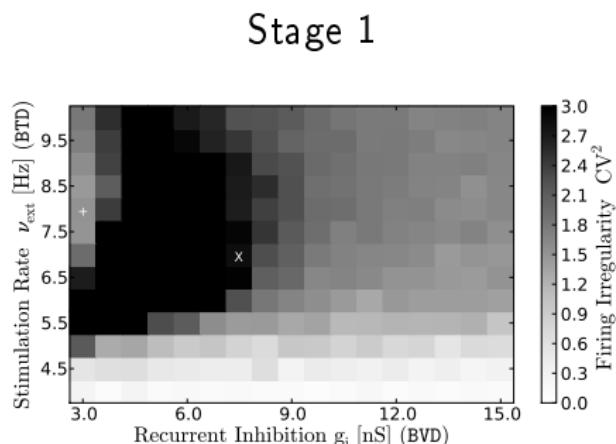
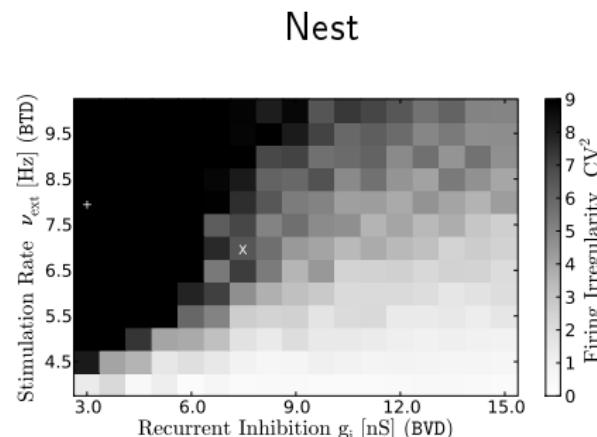
## Recurrent network dynamics

## Firing rate



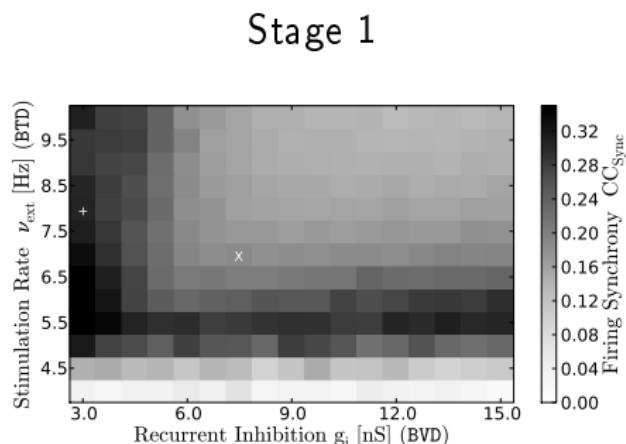
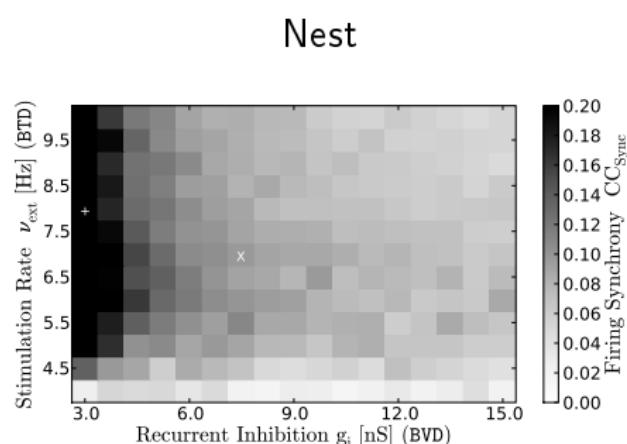
## Recurrent network dynamics

## Irregularity

(c) Hardware: Irregularity  $CV^2$ (d) NEST: Irregularity  $CV^2$

## Recurrent network dynamics

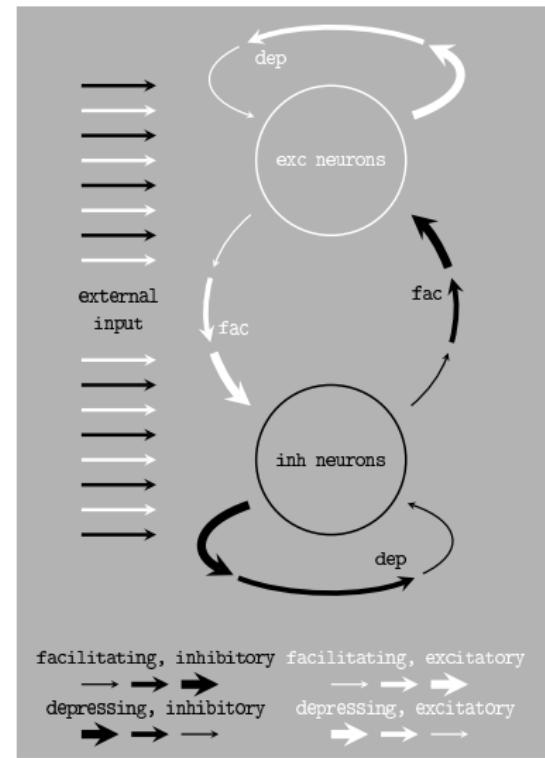
## Synchrony

(e) Hardware: Synchrony  $CC_{Sync}$ (f) NEST: Synchrony  $CC_{Sync}$

## Self-Stabilizing Network Architectures

# *Self-Stabilizing Network Architectures*

- Johannes Bill<sup>a</sup> & Klaus Schuch<sup>b</sup>
- inspired by Sussillo et al., 2007
- excitatory population  $P_e$  comprising  $N_e = 144$  neurons
- inhibitory population  $P_i$  consisting of  $N_i = 48$  neurons
- interconnected via depressing and facilitating synapses – see schematic

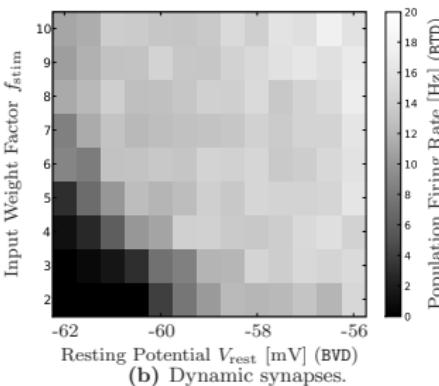
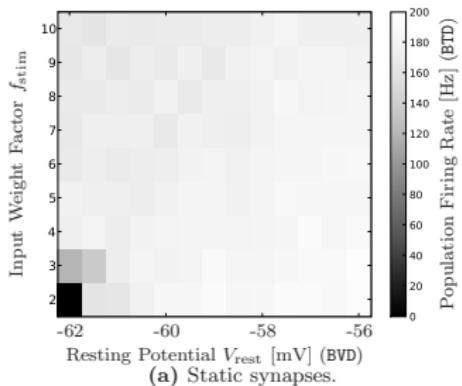


<sup>a</sup>University of Heidelberg/TU Graz

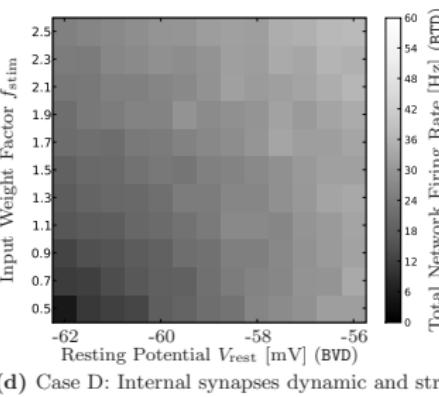
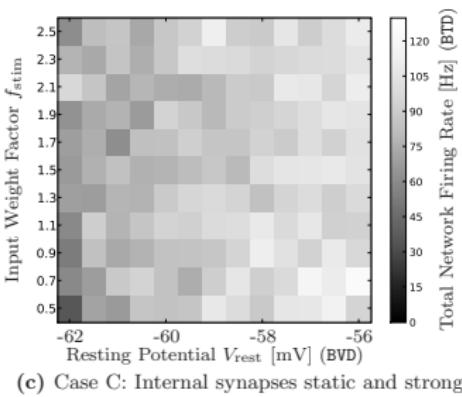
<sup>b</sup>TU Graz

## Self-Stabilizing Network Architectures

HW



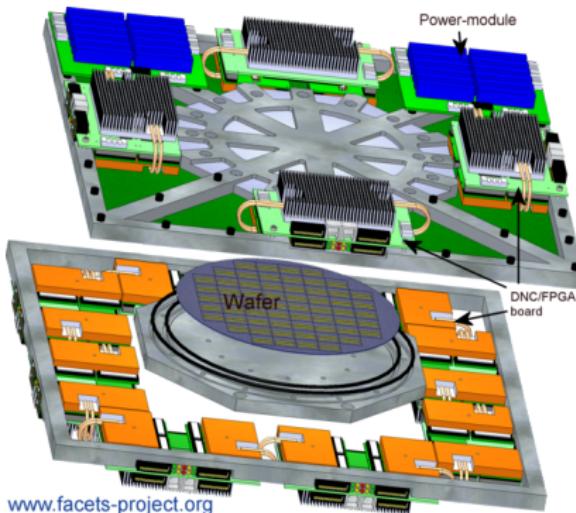
PC-SIM



# Goals

- comparing simulators and real hardware implementation
  - chip specification, debugging
- cooperation with external users
  - non-expert usability ⇒ PyNN

# FACETS Hardware Stage 2



[www.facets-project.org](http://www.facets-project.org)

- first HICANN<sup>a</sup> chip prototype produced
  - ~ 110 k hw-synapses,  
 $\leq 512$  AdEx neurons
  - $\#_{\text{inputs}} \times \#_{\text{neurons}} = \#_{\text{hw-synapses}}$
  - long- and short-term plasticity
  - accelerated operation (speed-up  $10^4$ )
- wafer system (2010?): ~ 350 HICANNS
- basic functionality test in preparation
- development of low level control software

---

<sup>a</sup>High Input Count Analog Neural Network

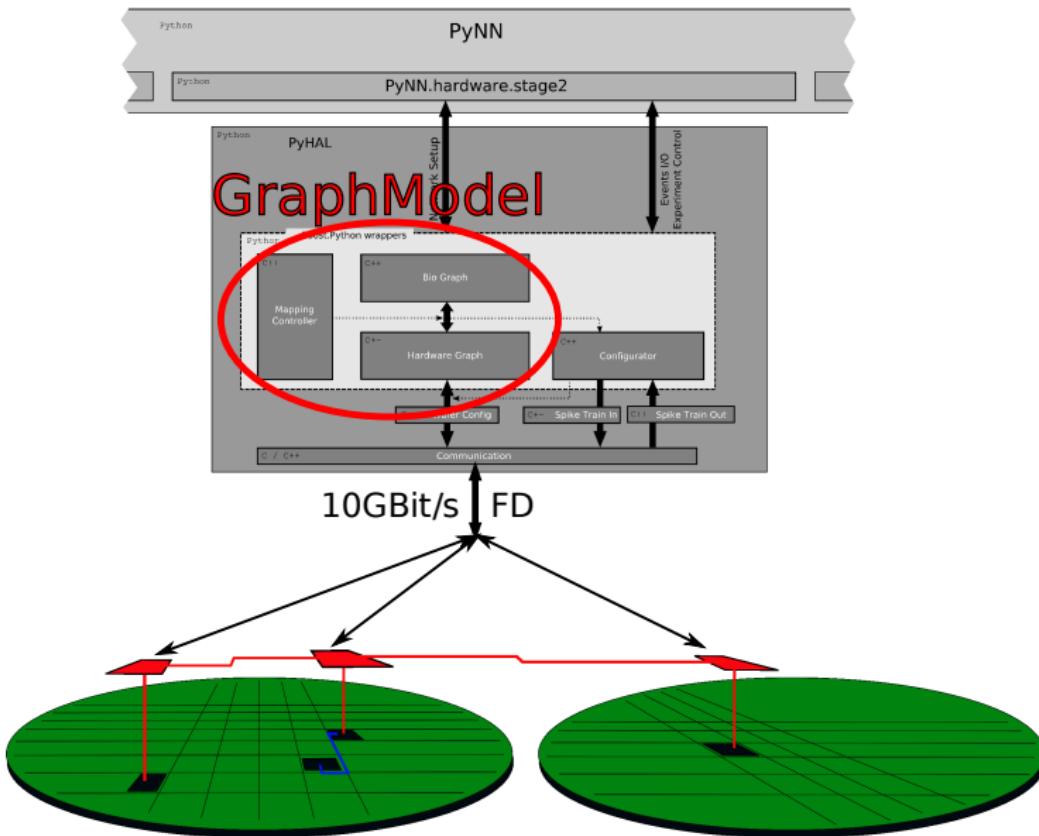
## Status – `pyNN.hardware.stage2`

- first use case: FACETS Demonstrator
  - an emulation of the real hardware system
  - low-level software connecting the hardware system and the existing simulation software stack is still missing
- procedural PyNN-API supported, except for
  - `record_gsyn`, `*_gsyn` – hardware constraint,  $g_{syn}$  not observable
  - `record_v` – hardware solution unknown (ADCs?) so far, recording via oscilloscope possible
- basic translation from object-oriented API to low-level API working

## Setup – Graph Model

Operating Interface

Mapping Software

Digital  
(Layer2)Analog  
(Layer1)

## Setup – Graph Model

## Bio model

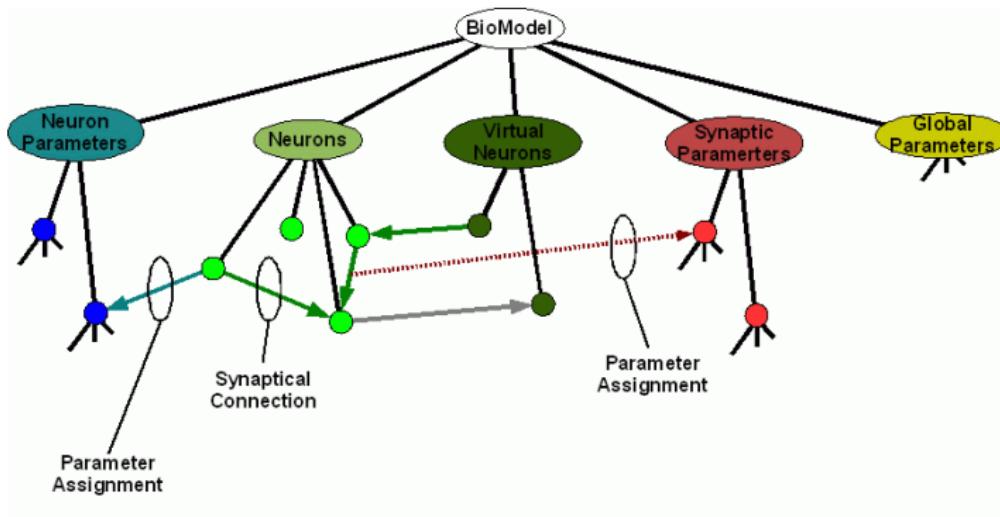


Figure by TUD

## Setup – Graph Model

## Hardware model

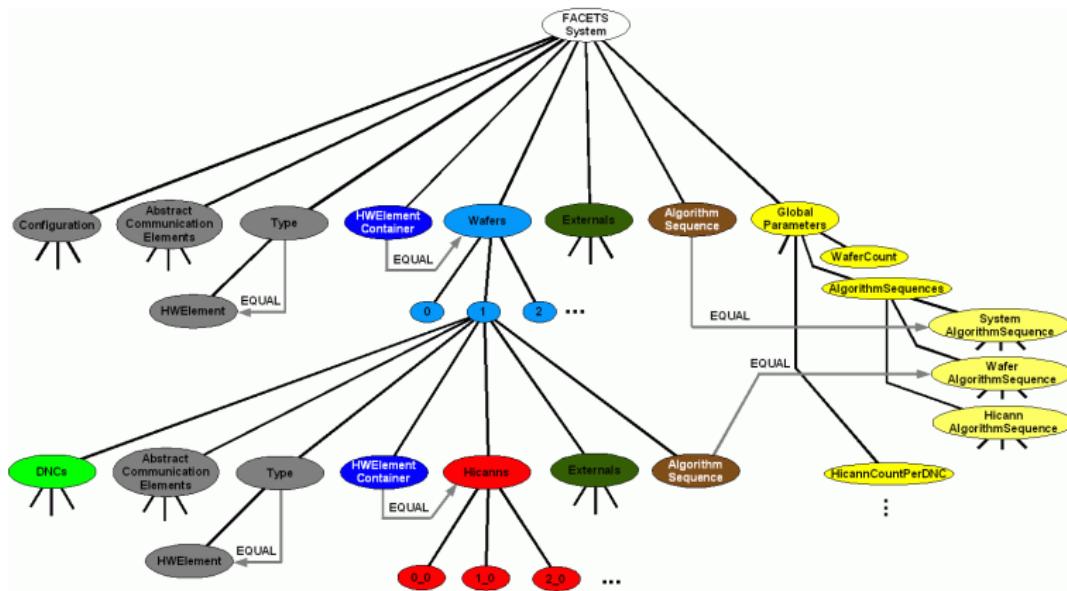
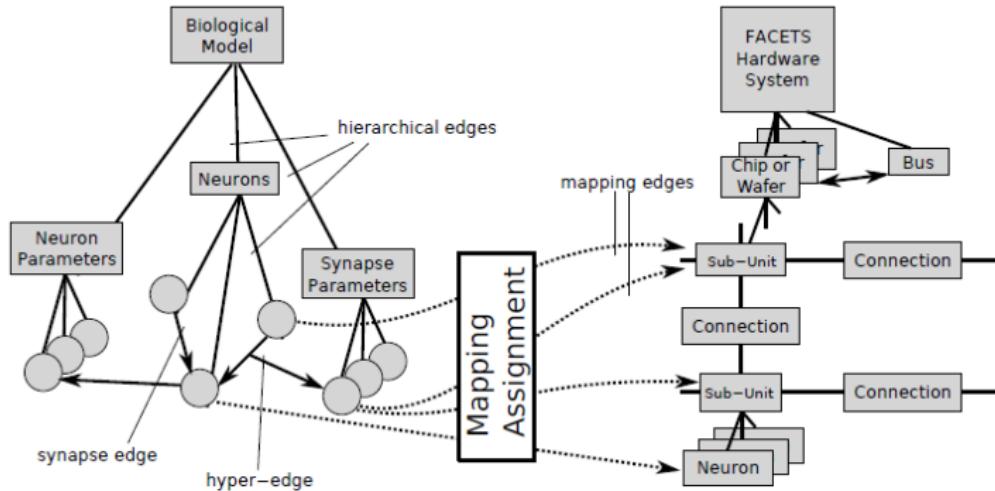


Figure by TUD

## Setup – Graph Model

## Mapping between models



**Setup – Routing**

# Routing

- 64 horizontal busses per HICANN row ⇒ approx. 1 k per wafer
- 256 vertical busses per HICANN column ⇒ approx. 9 k per wafer
- sparse switch matrices connect vertical, horizontal busses & HICANNs
- complex routing

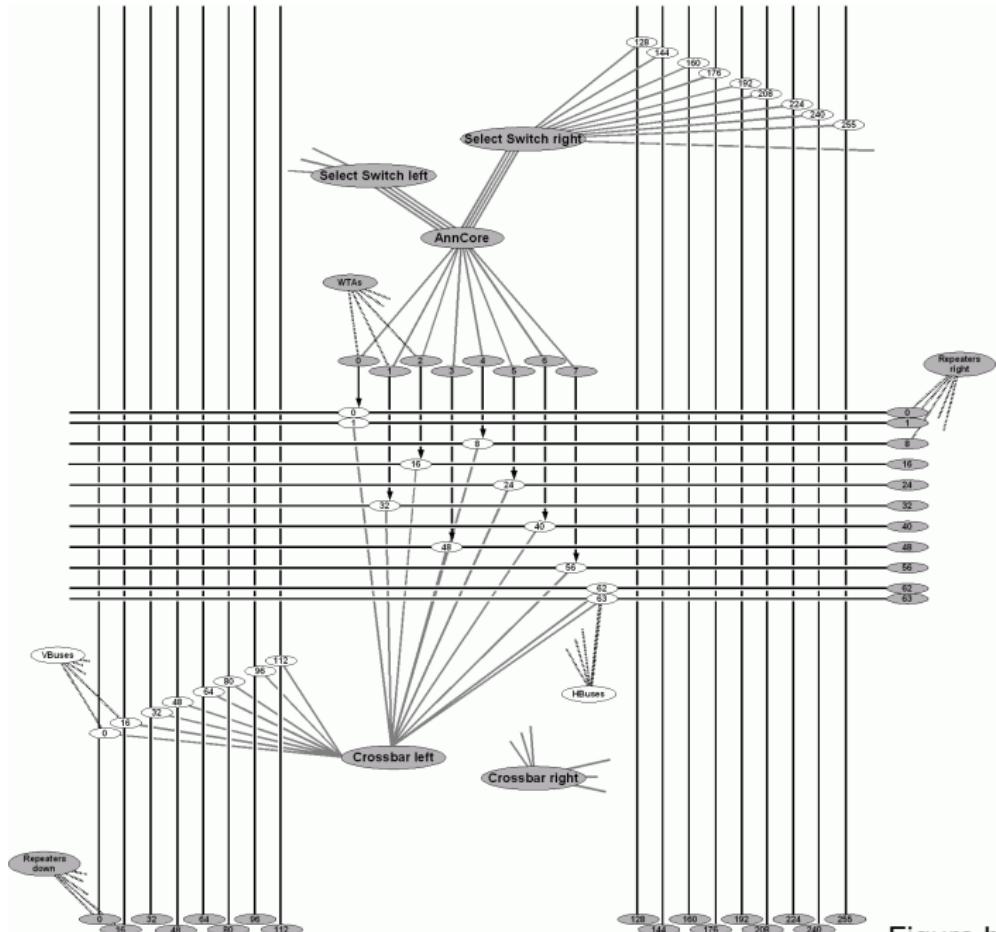
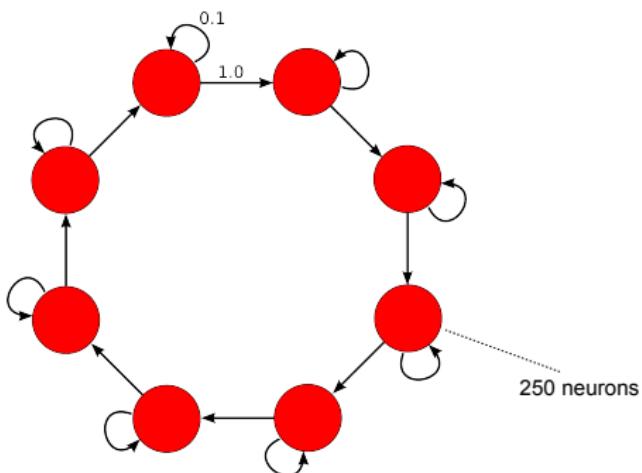


Figure by TUD

## Testing

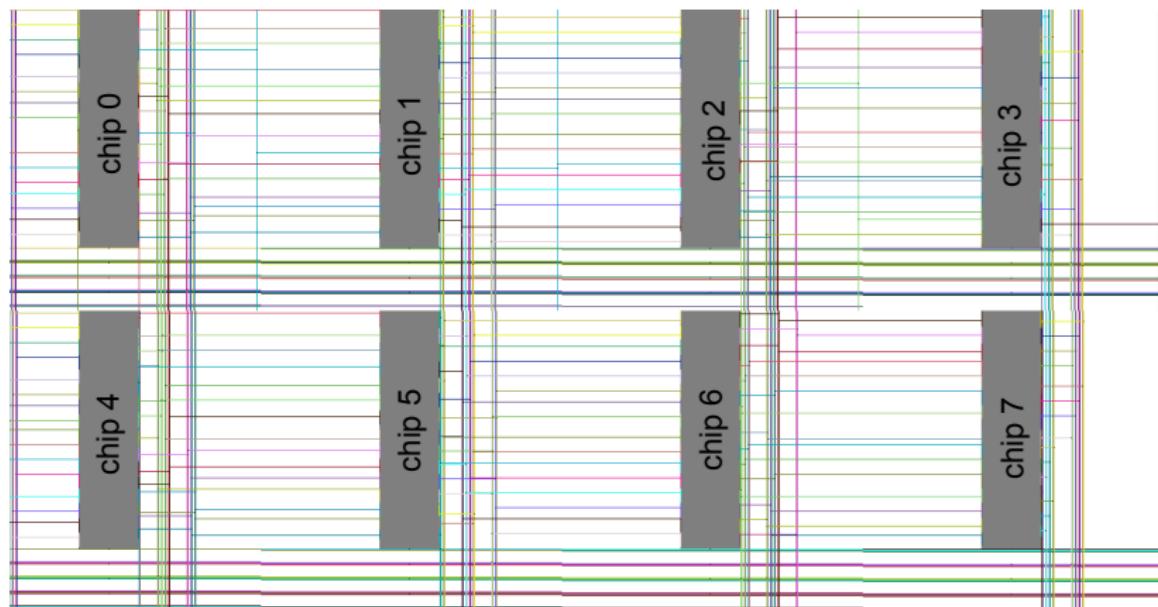
## A test of the placing and mapping steps



- dummy network example
- populations connected in a ring structure (e.g. synfire chain)

## Testing

## Routing after random placing



**Testing**

# *intelligent* mapping before routing

## NForceCluster Algorithm

- clusters neurons according to similarity criterions
- e.g. identical input sources

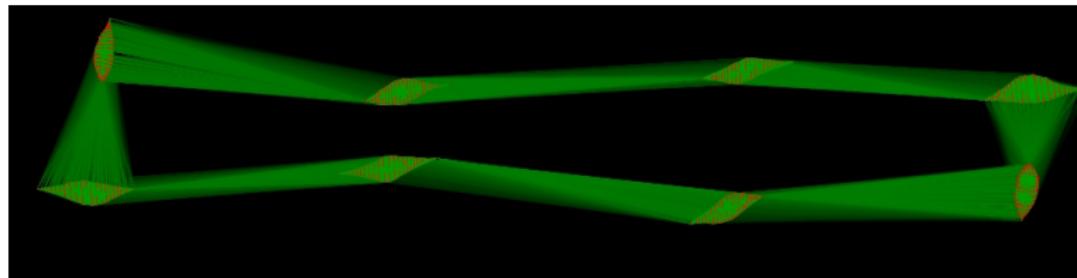
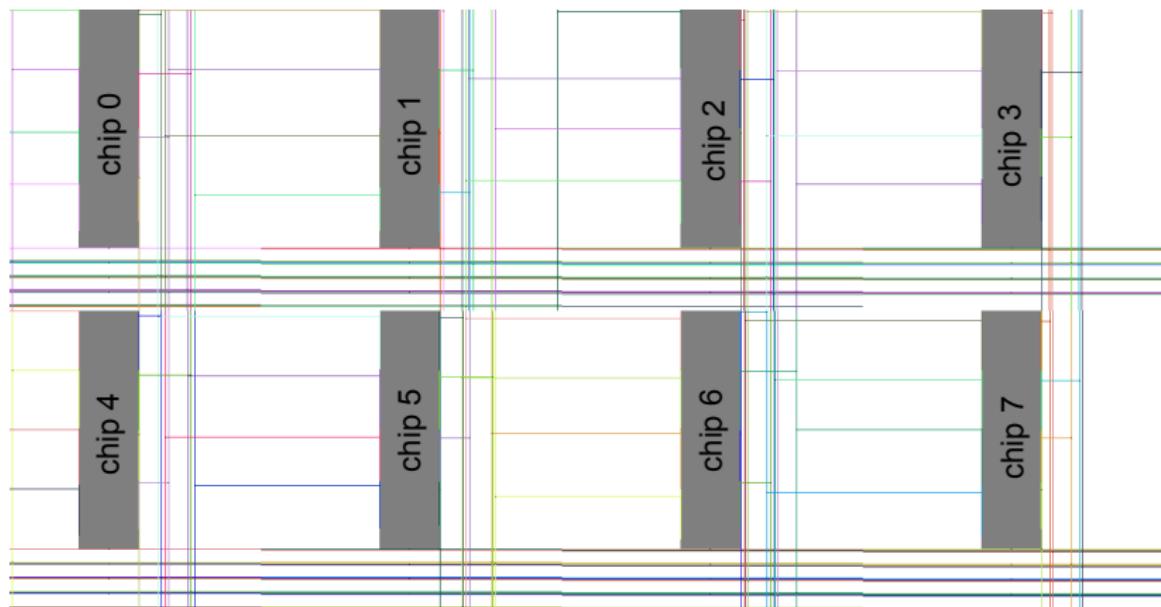


Figure by TUD

## Testing

routing after *intelligent* placing – NForceCluster Algorithm



## Testing

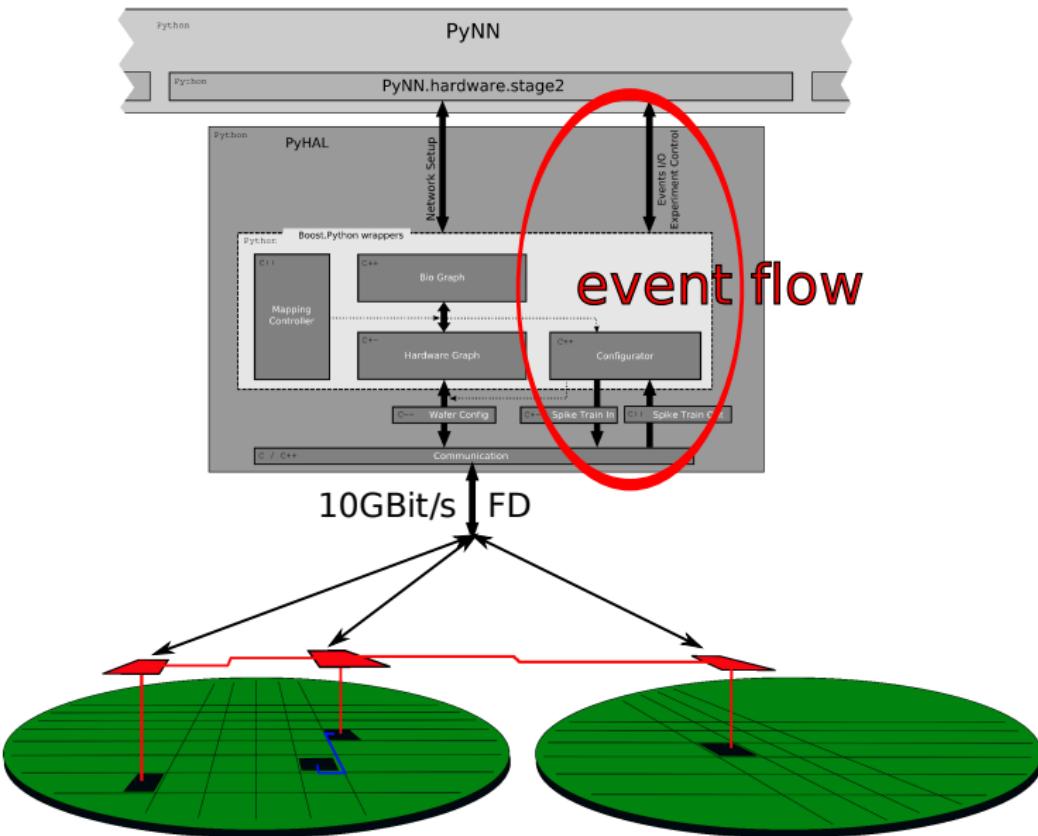
## Pre- &amp; postmap analysis

- debugging mapping and routing
  - measure for mapping quality
- 1 create representation of bio model in a flat format (PyNN script, procedural API, every connection, every parameter in a single row)
  - 2 run mapping and routing steps (synaptic loss due to connectivity constrains)
  - 3 create representation of the mapped hardware model
  - 4 differences between scripts due to
    - neuron loss, synaptic loss (limited bandwidth and routing resources)
    - parameter distortion (limited parameter range, flexibility, delays)

## Runtime control – event flow

Operating Interface

Mapping Software

Digital  
(Layer2)Analog  
(Layer1)

## Runtime control – event flow

- after setup, the *Configurator* controls event flow
  - input- and output spikes
  - runtime chip control (recording, read-out of parameters, setting parameters)

## Incremental configuration in PyNN

# Incremental configuration in PyNN

- main problem: clustering of parameter space takes a lot of time
- placing and routing steps are slow (1 min – 1 h, depends on model)
- but parameter changes can be typically done at zero cost ⇒ parameter sweeps still possible
- some topological changes do not require rerunning the mapping steps
- we will need support for *incremental* configuration: in our interface: `setup()`, `connect()`, `run()`, `reconnect()`, `run()`, ...
- PyNN support for `disconnect()`, setting and resetting parameters of high-level objects?

## Backend-specific pre-processing of data

# Backend-specific pre-processing of data

- multi-wafer systems will generate too many spikes to analyze (or store) in real-time
- even a single wafer can't be fully recorded
  - on-wafer bandwidth is in the range of several TB/s
  - off-wafer (FPGA) bandwidth: ~ 50 GB/s
  - host-FPGA bandwidth: 2 GB/s
- reduce data rates – FPGA can do a lot of preprocessing
- statistical and averaged measures evaluated online
  - avg. firing rates, isi, measures for synchrony, irreg., etc.
  - low-pass filtering of analog (e.g. sub-threshold membrane trace) data, avg. weight over time

pyNN.hardware.stage2 – Open issues

## Conclusion: Open issues

- hardware specific unit tests
- fragile make flow needs to be replaced by a better build system
  - we have to link to the GraphModel:  
C++ and dependencies on other libraries
- create dummy PyNN module for virtual hardware (using the FACETS Demonstrator simulation code)
- lots of control software missing
- waiting for the real wafer system