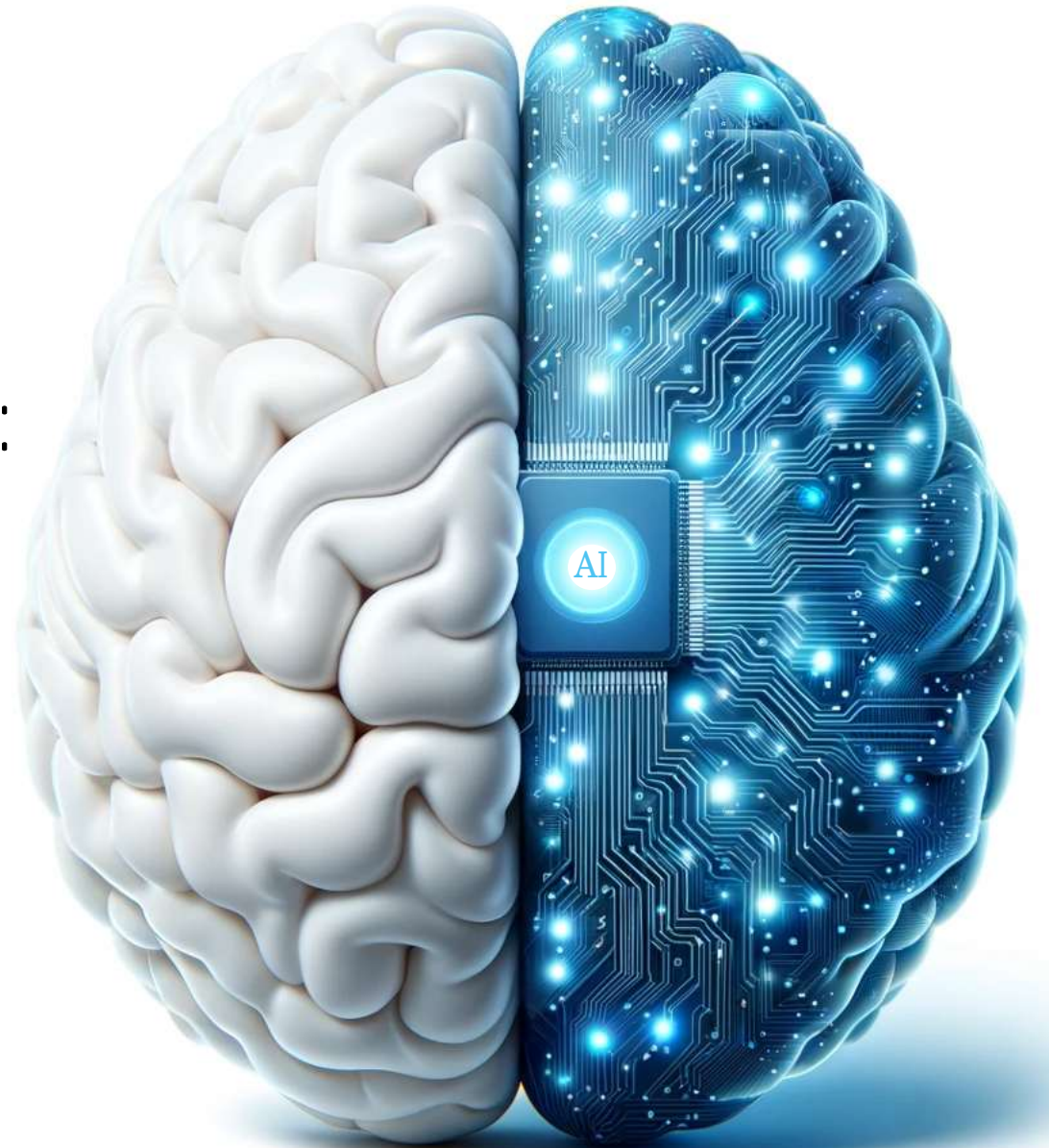


SpiNNcloud

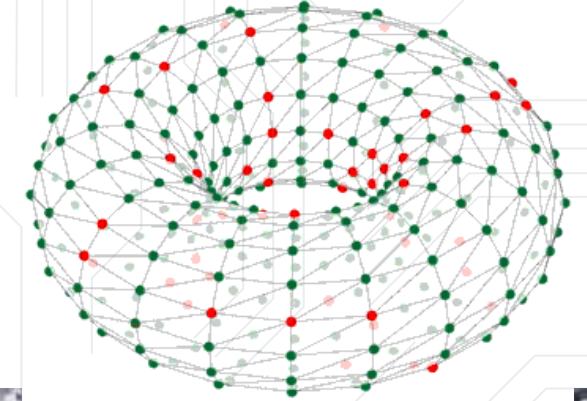
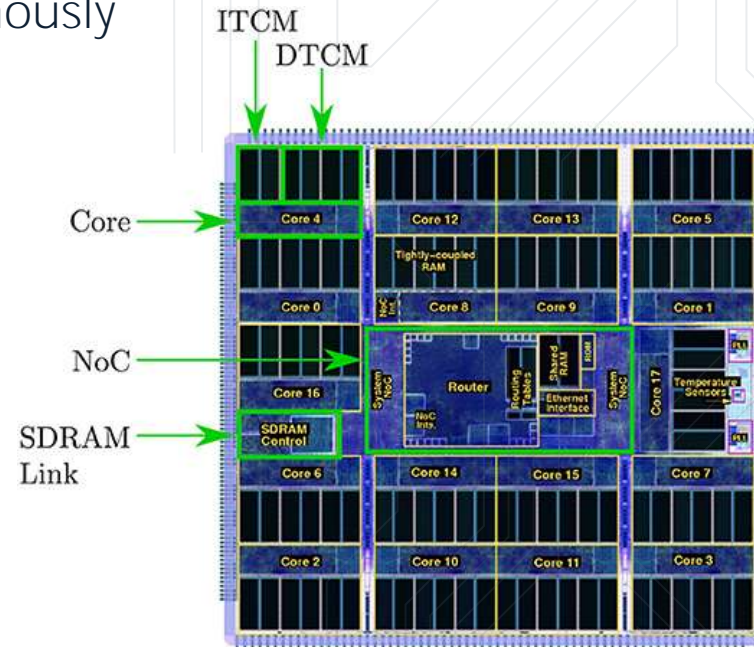
Brain-inspired computing:
Systems for the
next generation of AI



Matthias Lohrmann
CTO

Redefining Brain-inspiration: SpiNNaker (1)

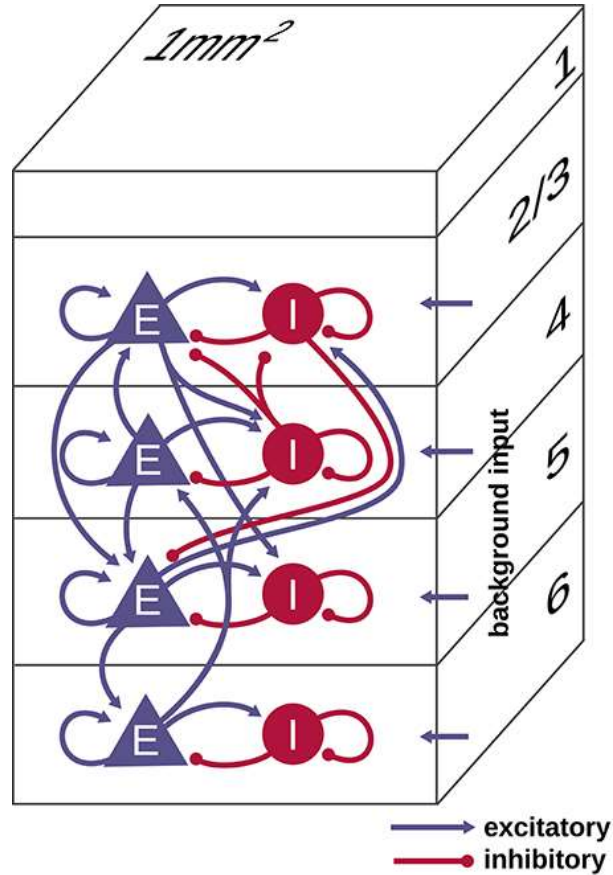
- Do spiking neurons digitally & asynchronously
- 18 ARM 968 processors / chip
- 1 router / chip with world model
- 1M core supercomputer built



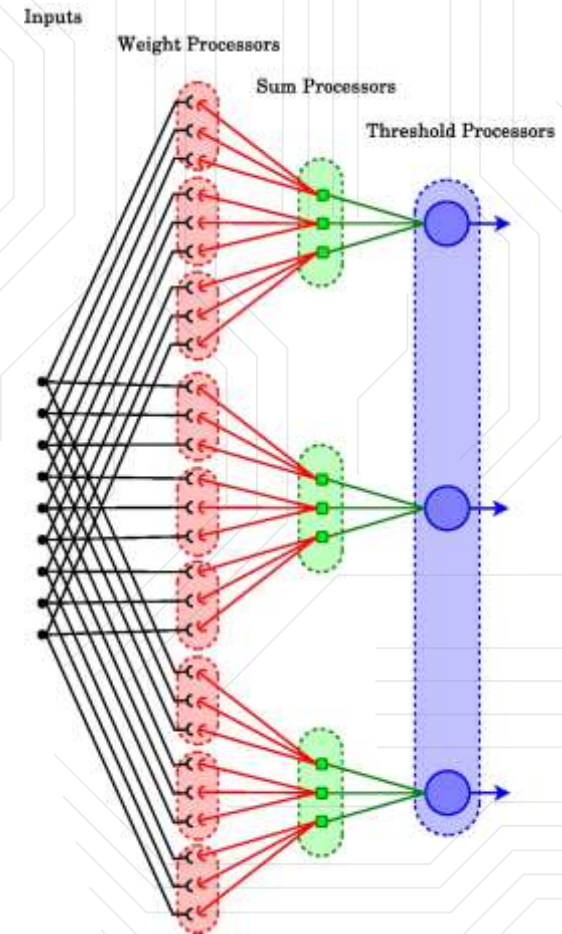
<http://apt.cs.manchester.ac.uk/projects/SpiNNaker/SpiNNchip/>
<http://apt.cs.manchester.ac.uk/projects/SpiNNaker/project/Access/>



SpiNNaker 1



- Cortical Brain Model



- Multilayer Perceptron

Learnings for AI

Granularity

- MIMD
- hybrid systems
- exploit sparsity

Sparsity

- activation
- weights
- representation

Functionality

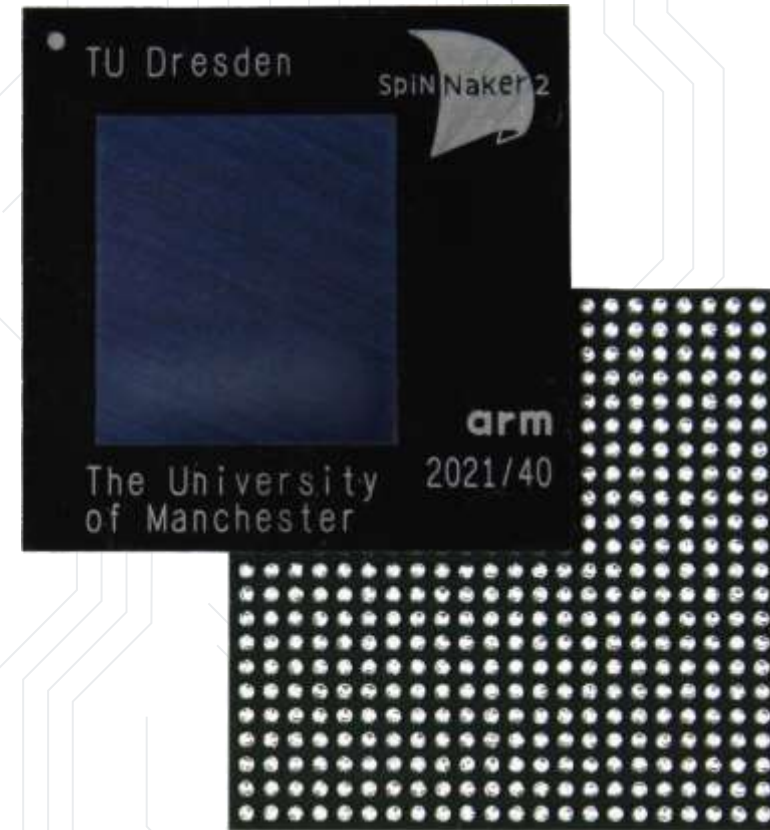
- memory dominates area
- compute is for free

Unique Hybrid Microchip SpiNNaker2

- Faster than NVIDIA's A100 in brain models
- Consumes 1/10th less power than GPUs
- Enables more-than-DNN AI systems



IP is backed by multiple international patents

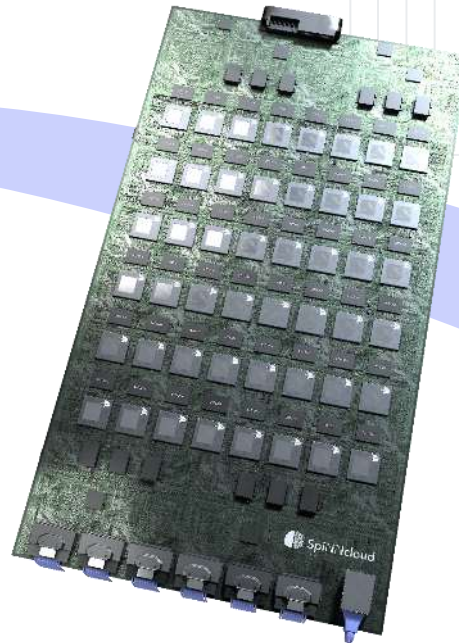


Unique Hybrid Microchip SpiNNaker2

SpiNNcloud



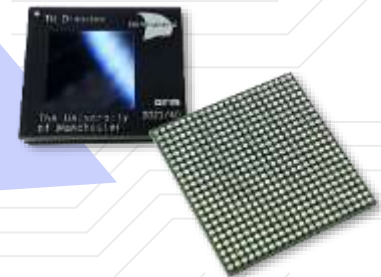
Cloud Board



SpiNNnode

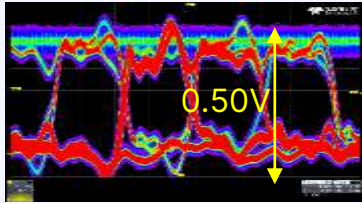


SpiNNaker2



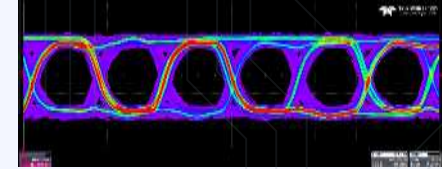
IP is backed by multiple international patents

The SpiNNaker2 Architecture



6x6 Serial Chip-2-Chip Links

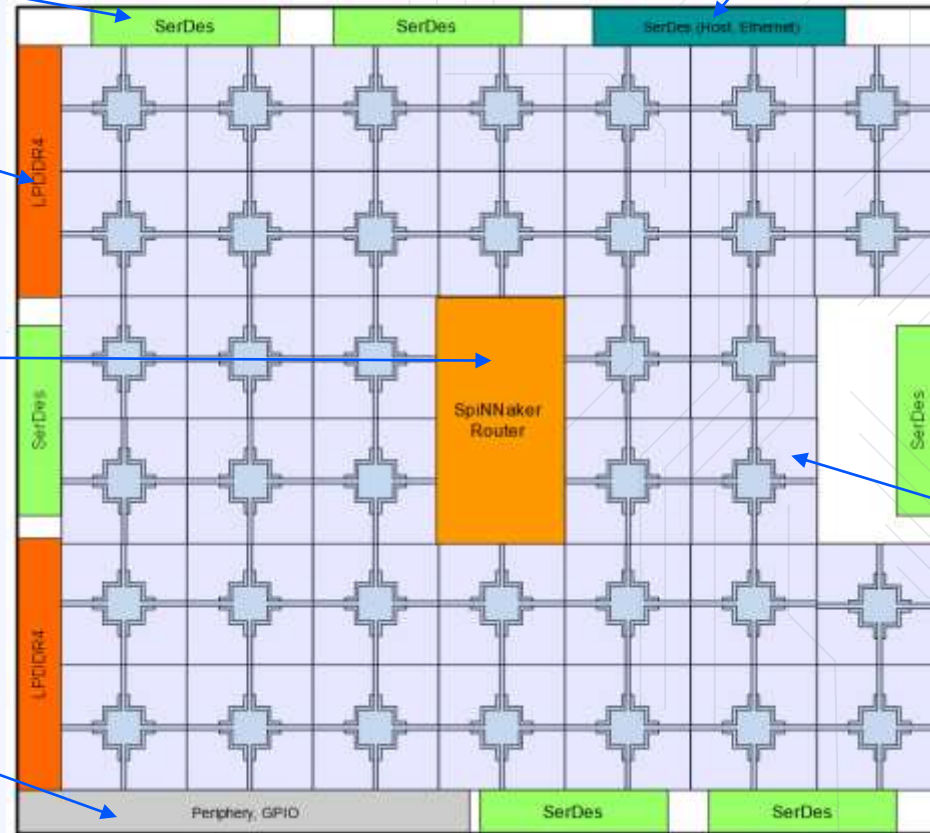
- SerDes
- Board2Board coms
 - Gbit Ethernet



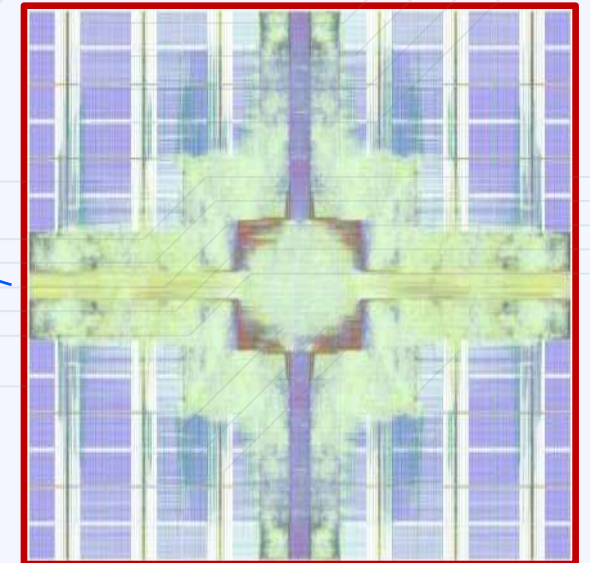
2x LPDDR4 Memory Interfaces

- SpiNNaker Packet Router
- Lightweight event packet communication fabric

- Periphery
- Flexible GPIO, QSPI, I2C (Master + Slave), JTAG



152 ARM M4F cores + accelerators



The SpiNNaker2 Architecture

Dynamic Power Management

- DVFS and PSO

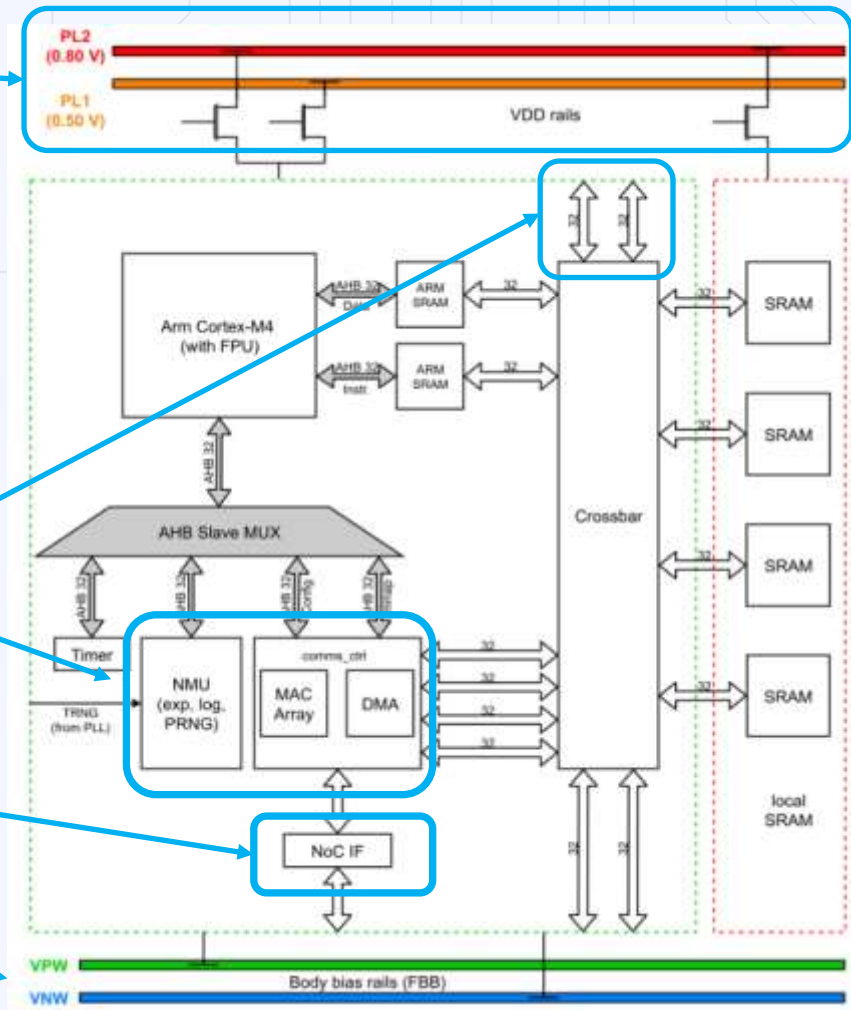
Accelerators

- MAC + DMA
- Exp / Log
- Random numbers (PRNG, TRNG)

Communication

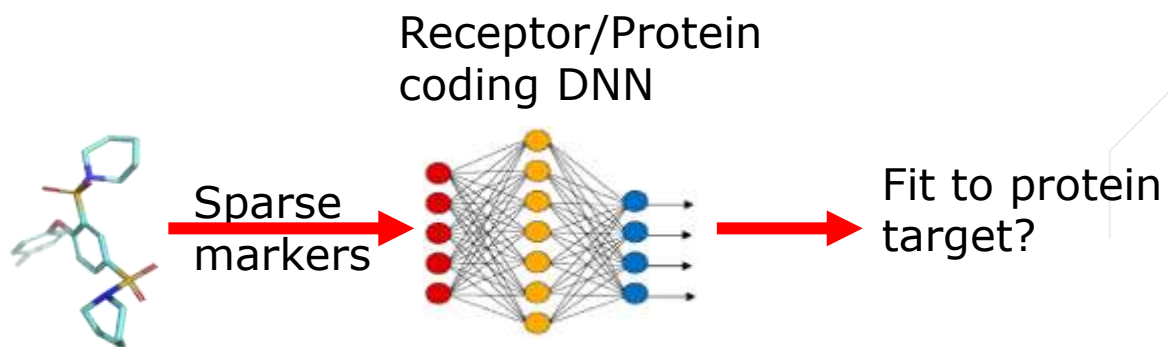
- Neighbour PE memory sharing (synchronous)
- On- and off-chip memory access
- Event & data handling

ABX Adaptive Body Biasing

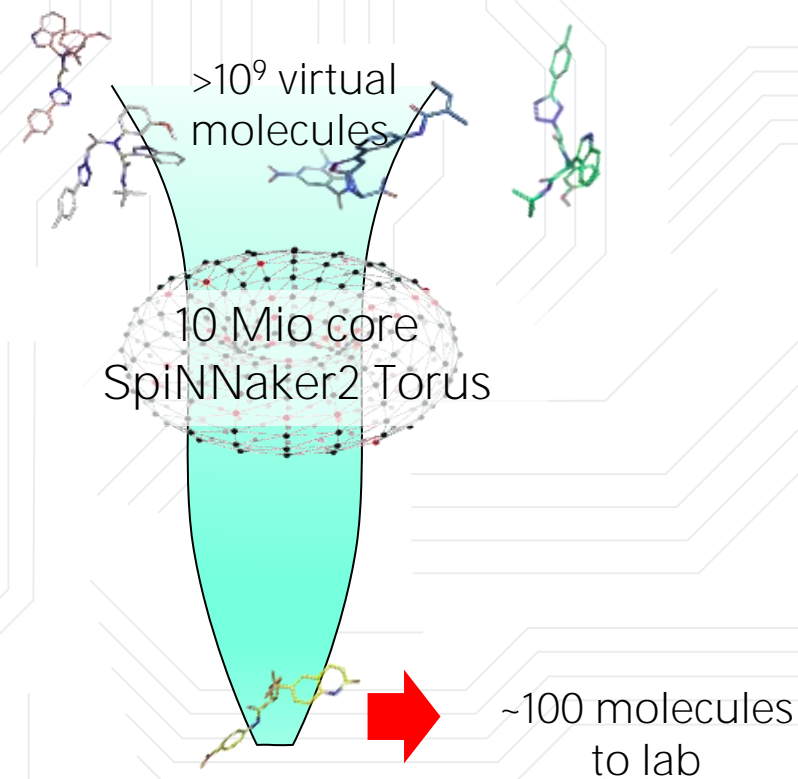


Hybrid AI/non-AI Compute: Drug Discovery

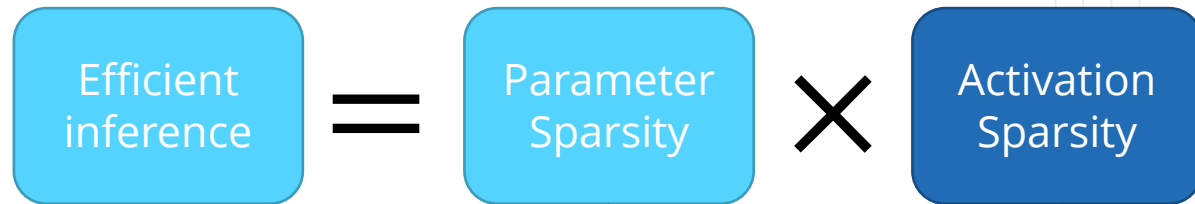
- Target: 10^9 molecule scanning in O(hours) instead of O(weeks)
- Densely couple CPU + GPU for speed up



Granularity !



Efficient inference of large language models

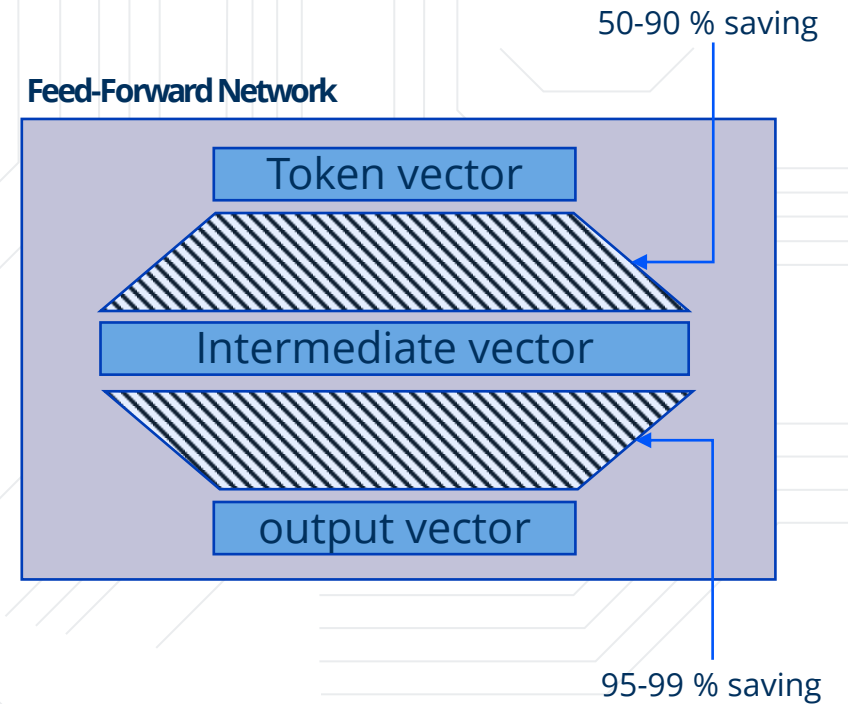


50 % savings (one-shot i.e. no retraining)

Frantar & Alistarh
"SparseGPT: Massive Language Models
Can be Accurately Pruned in One-Shot"

75-95 % savings

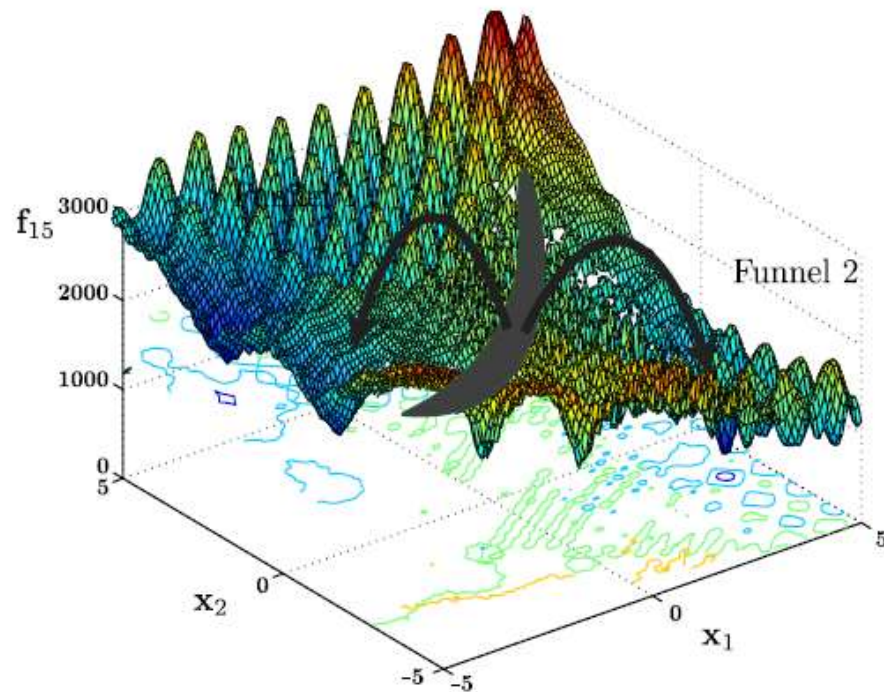
Zhu et al.
"SpikeGPT: Generative Pre-trained
Language Model with Spiking Neural Networks"



Sparsity !

Particle Swarm Evolutionary Algorithms

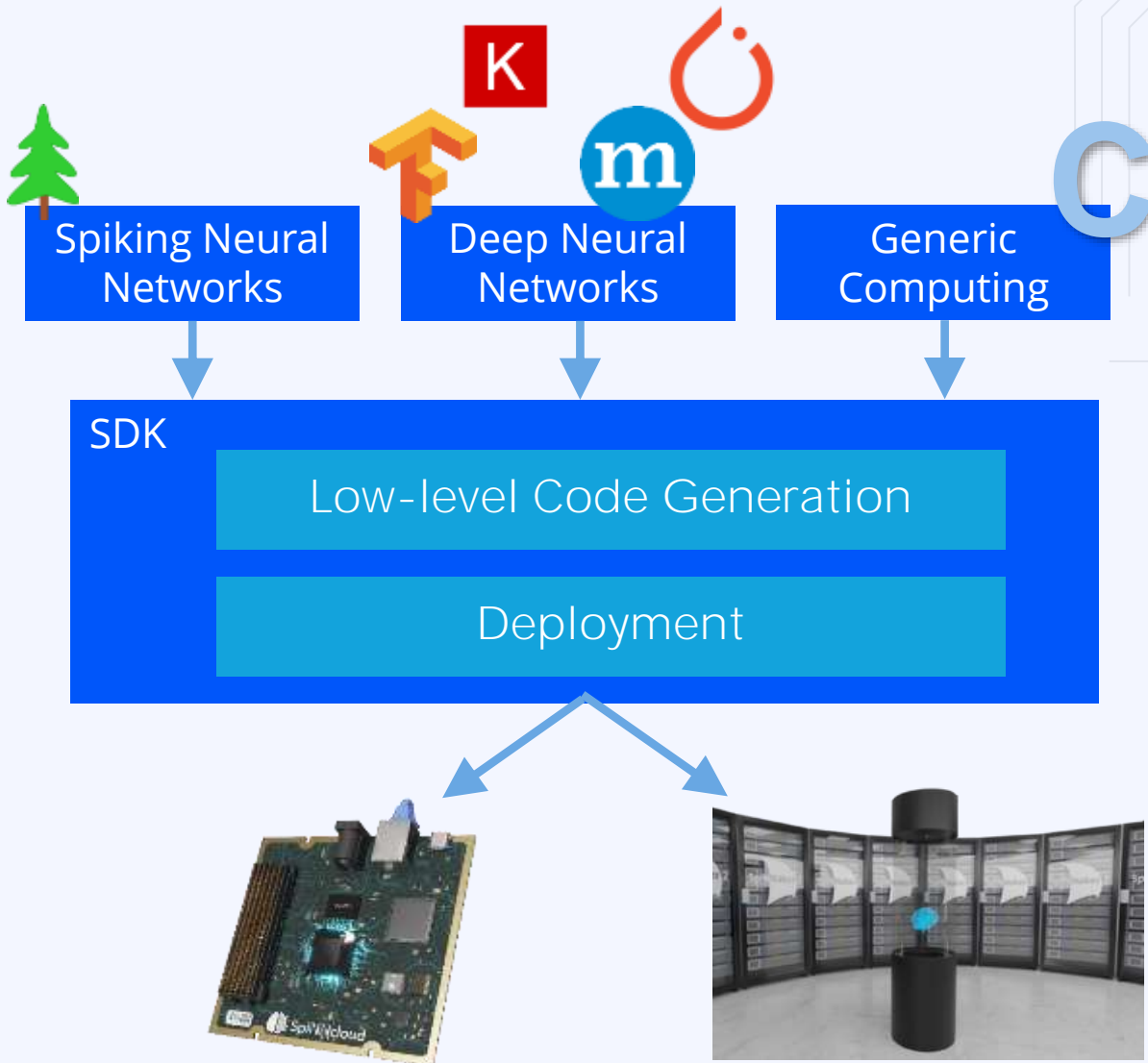
- Island-based, distributed genetic algorithms for optimization
- Combine the robust local search performance the global exploration power of PSO (particle-swarm optimization)



Functionality!

Fig. 1. Two-dimensional version of the highly dispersive function f_{15} from the CEC benchmark test suite [17]. The global topology is a double funnel separated by the central ridge region (in gray). The global and several local minima are contained in funnel 1, several deep local minima in funnel 2. This topology is hard since a search heuristic can be trapped in the broad funnel 2.

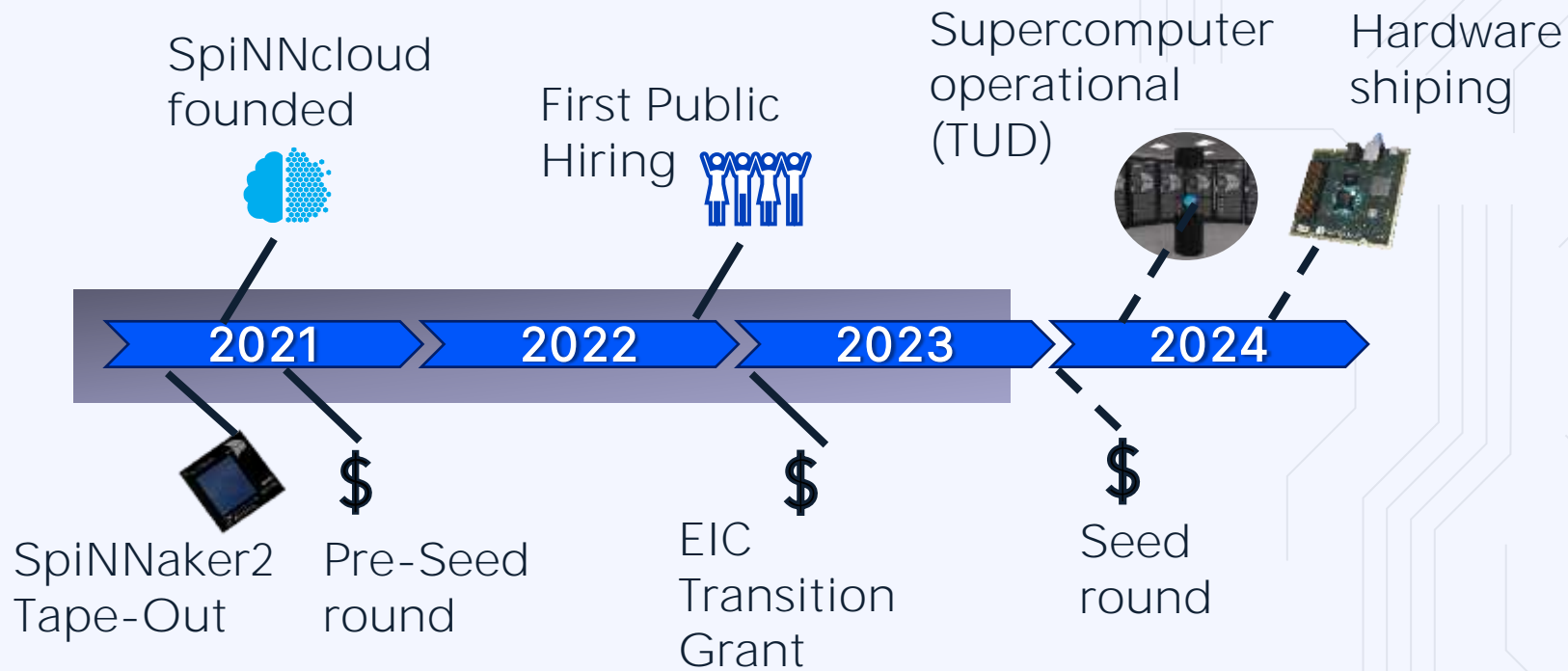
Unified Software Stack



- Simple development + deployment with high flexibility
- Full suite available in Summer 2024

```
1 from spinnaker2 import snn, hardware
2
3 neuron_params = {
4     "threshold":1.,
5     "alpha_decay":0.9,
6 }
7
8 stim = snn.Population(
9     size=10,
10    neuron_model="spike_list",
11    params={0:[1,2,3], 5:[20,30]},
12    name="stim")
13
14 pop1 = snn.Population(
15    size=20,
16    neuron_model="lif",
17    params=neuron_params,
18    name="pop1")
```

About SpiNNcloud



Looking for

- Follow investors for seed round
- Compiler engineering talent
- Proof-of-concept partners

Thank You!



Contact: matthias.lohrmann@spinncloud.com

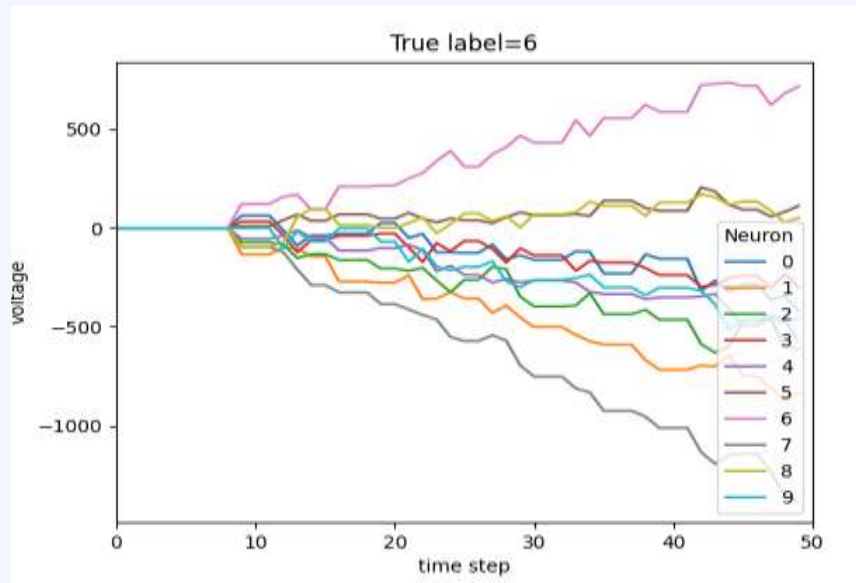
SpinNcloud



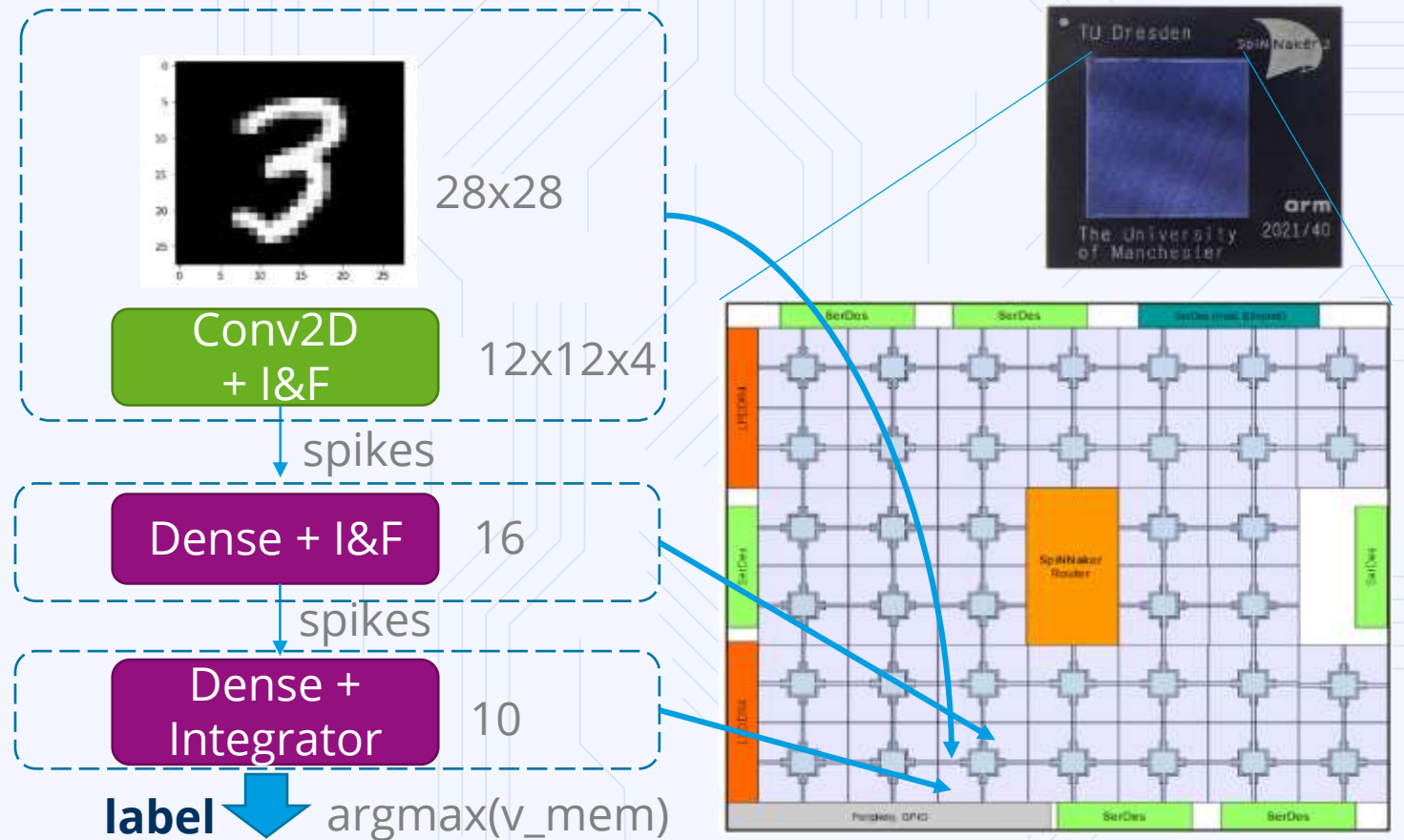
Unified Software Stack

conv2d if neuron rate

lif neuron



Accuracy identical to DNN



Particle Swarm Evolutionary Algorithms

- Island-based, distributed genetic algorithms for optimization
- Combine the robust local search performance the global exploration power of PSO (particle-swarm optimization)

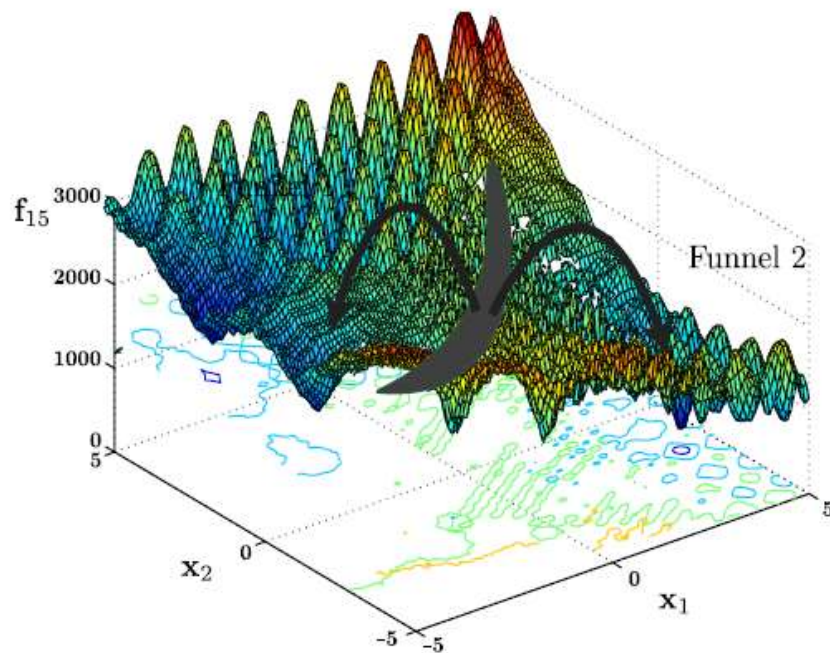
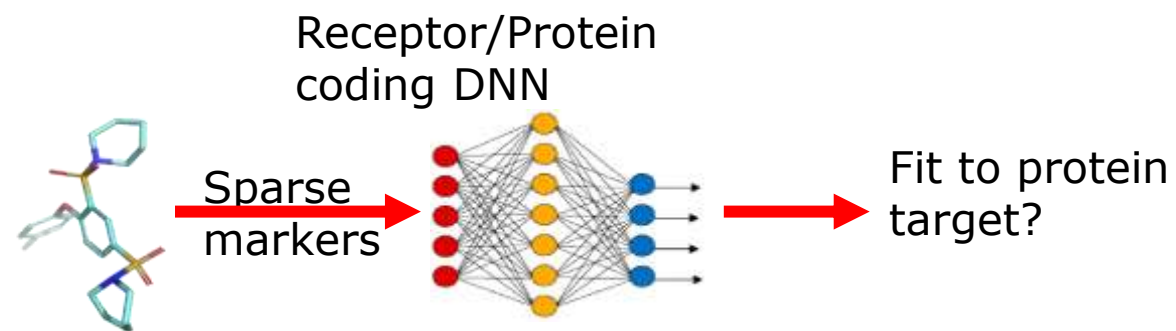
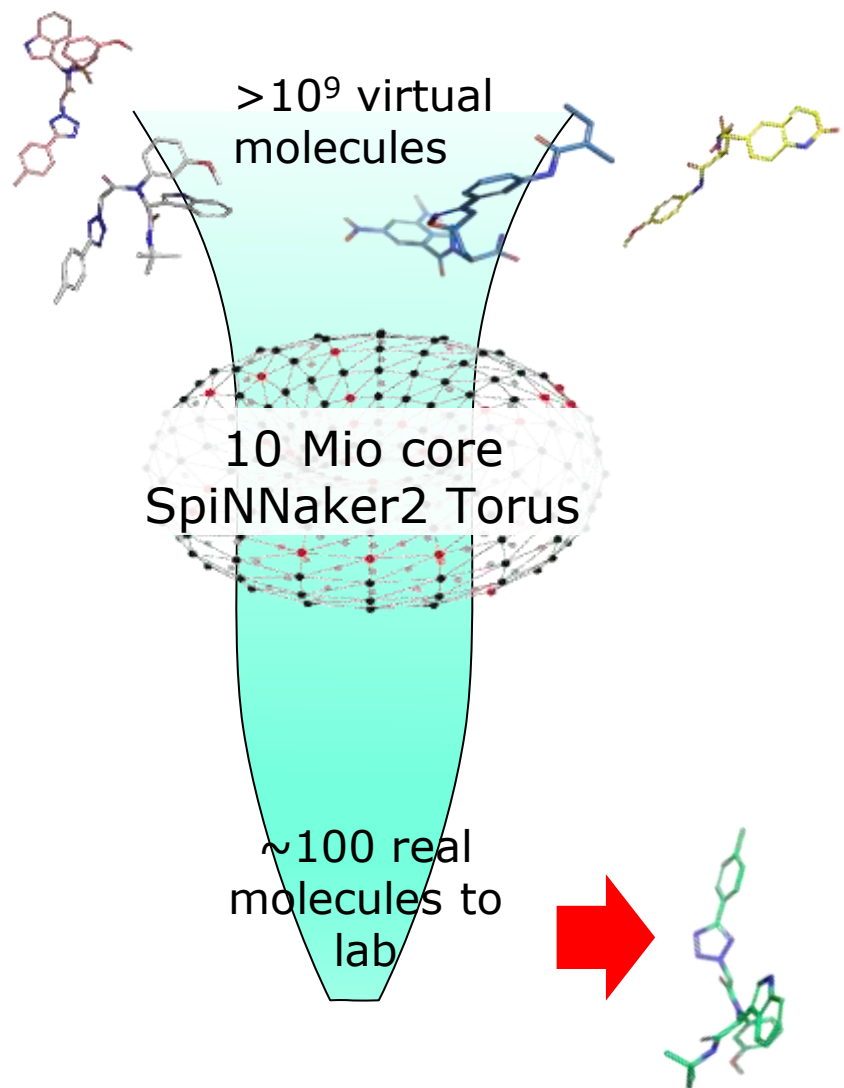


Fig. 1. Two-dimensional version of the highly dispersive function f_{15} from the CEC benchmark test suite [17]. The global topology is a double funnel separated by the central ridge region (in gray). The global and several local minima are contained in funnel 1, several deep local minima in funnel 2. This topology is hard since a search heuristic can be trapped in the broad funnel 2.



Del Alamo, Sala, McHaourab, **Meiler**: "Sampling alternative conformational states of transporters and receptors with AlphaFold2"; *Elife*; **2022** || Brown, Mendenhall, Geanes, **Meiler**; "General Purpose Structure-Based Drug Discovery Neural Network Score Functions with Human-Interpretable Pharmacophore Maps"; *J Chem Inf Model*; **2021** ||