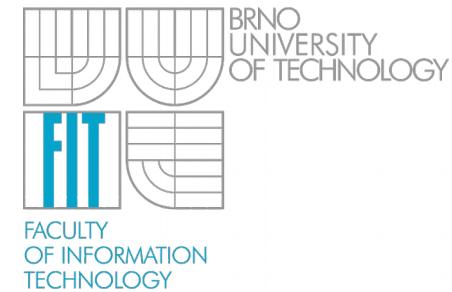


A formal model of FPGA implemented artificial neural networks

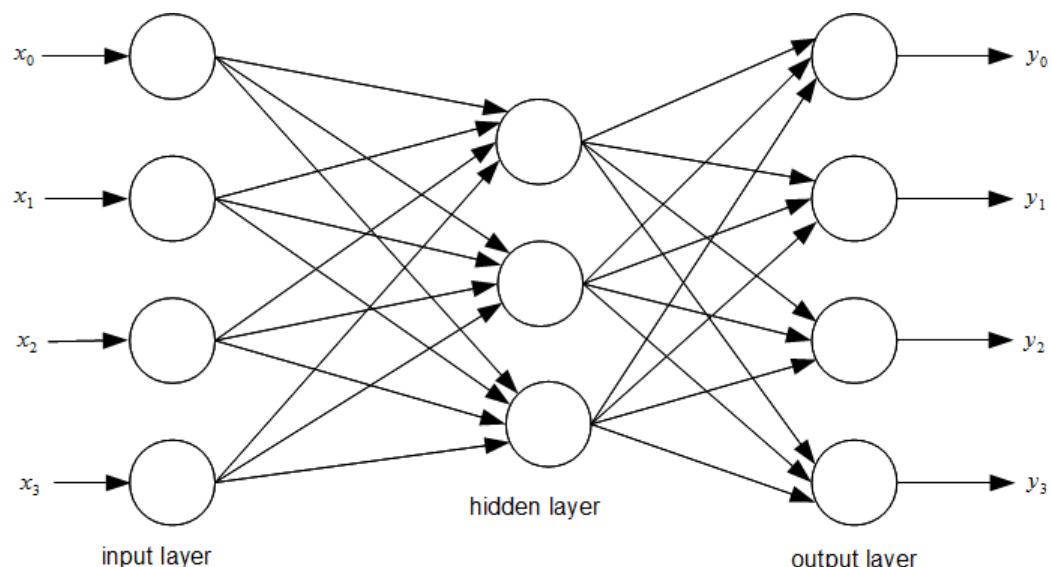
Brno University of Technology
Faculty of Information Technology
Božetěchova 2, 612 66 Brno
Martin Krčma, ikrcma@fit.vutbr.cz



Outline

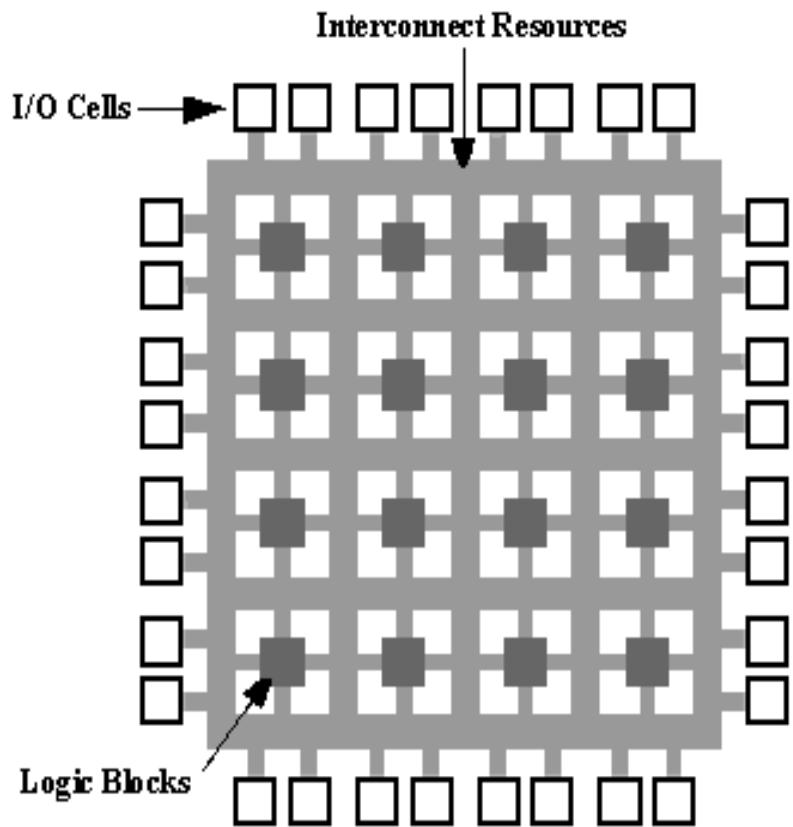
- **Introduction**
 - Introduction of artificial neural networks.
 - Introduction of FPGAs.
- **Neural networks in FPGAs**
 - Definition of FPNA.
 - Definition of FPNN.
 - Definition of grid FPNN.
 - Grid FPNN mapping algorithm for redundancy free fault tolerance.
- **Summary**

- Computational model of the human brain
- History
 - introduced in 1943 (Warren McCulloch, Walter Pitts)
 - first conference in 1987
- Usage
 - classification
 - prediction
 - approximation
 - control
 - clustering



- Field programmable gate array.
- Digital circuits repeatedly programmable on low level (hardware level).
- Capable of implementing almost every type of digital circuits.
 - CPU, GPU, DSPs, filters, controllers...
- Capable of parallel computing.
- Much faster than general purpose processors in parallel task designs.
- Not as fast as ASIC, but flexible, repetitively usable.
- Great development and experimental tool.

Introduction: FPGA



- Field programmable neural array.
- Introduced by Bernhard Girau in his master thesis.
- FPGA implementation of neural networks.
- Can be also a resources saving approximation.
- Relaxed definition.
- Strongly optional structure.

- Let N be a set of nodes (**activators**) and E a set of oriented edges (**links**). We say that graph (N, E) is an **FPNA** if the following statements hold:
 - Every node has a set of predecessors:
$$\forall n \in N : \exists \text{Pred}(n) = \{p \in N | (p, n) \in E\}$$
 - Every node has a set of successors:
$$\forall n \in N : \exists \text{Succ}(n) = \{s \in N | (n, s) \in E\}$$

Notes:

- The $e \in N \cup E$ is called **neural resource**
- The links are named as (x, y) where $x, y \in N$

- There is a set of input nodes:

$$\exists N_i = \{n \in N \mid \text{Pred}(n) = \emptyset\}$$

- Every link has affine operators:

$$\forall (r, n) \in E \wedge p \in \text{Pred}(n) : \exists \alpha_{(p, n)} = W_n(p) \cdot x + T_n(p)$$

- Every non-input activator has an iterative operator for potential computing:

$$\forall n \in N : \exists i_n : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$$

- Every non-input activator has a function operator for activation function computing:

$$\forall n \in N : \exists f_n : \mathbb{R} \rightarrow \mathbb{R}$$

| FPNA: it is not enough

- The definition of FPNA is relaxed and partially unspecific.
- It defines the existence of the activators and links operators.
- But it **does not** define their concrete values.
- Using *Pred* and *Succ* sets it defines the order of neural resources but it **does not** define the structure of their interconnection.
- It **does not** define the default values.
- Thus for creating a fully specified network (design) there is needed some more specification than FPNA offers.
- The remaining details specify the Field programmable neural Network – FPNN.

- Let N be a set of nodes (**activators**) and E a set of oriented edges (**links**). We say that graph (N, E) is an **FPNN** if the following statements hold:

- (N, E) is an FPNA
- default value of activators iteration variable:

$$\forall n \in N : \exists \theta_n \in \mathbb{R}$$

- a number of iterations:

$$\forall n \in N : \exists a_n \in \mathbb{N}$$

- concrete values of: $\forall \alpha_{(p, n)} : W_n(p) \in \mathbb{R} \wedge T_n(p) \in \mathbb{R}$
- a number of inputs for every input node:

$$\forall n \in N_i : \exists c \in \mathbb{N}$$

- binary flag determining interconnection between an activator n and link (p,n) :

$$\forall p, n \in N \wedge p \in \text{Pred}(n) : \exists r_n(p) \in \{0,1\}$$

- binary flag determining interconnection between an activator n and link (n,s) :

$$\forall n, s \in N \wedge s \in \text{Succ}(n) : \exists s_n(s) \in \{0,1\}$$

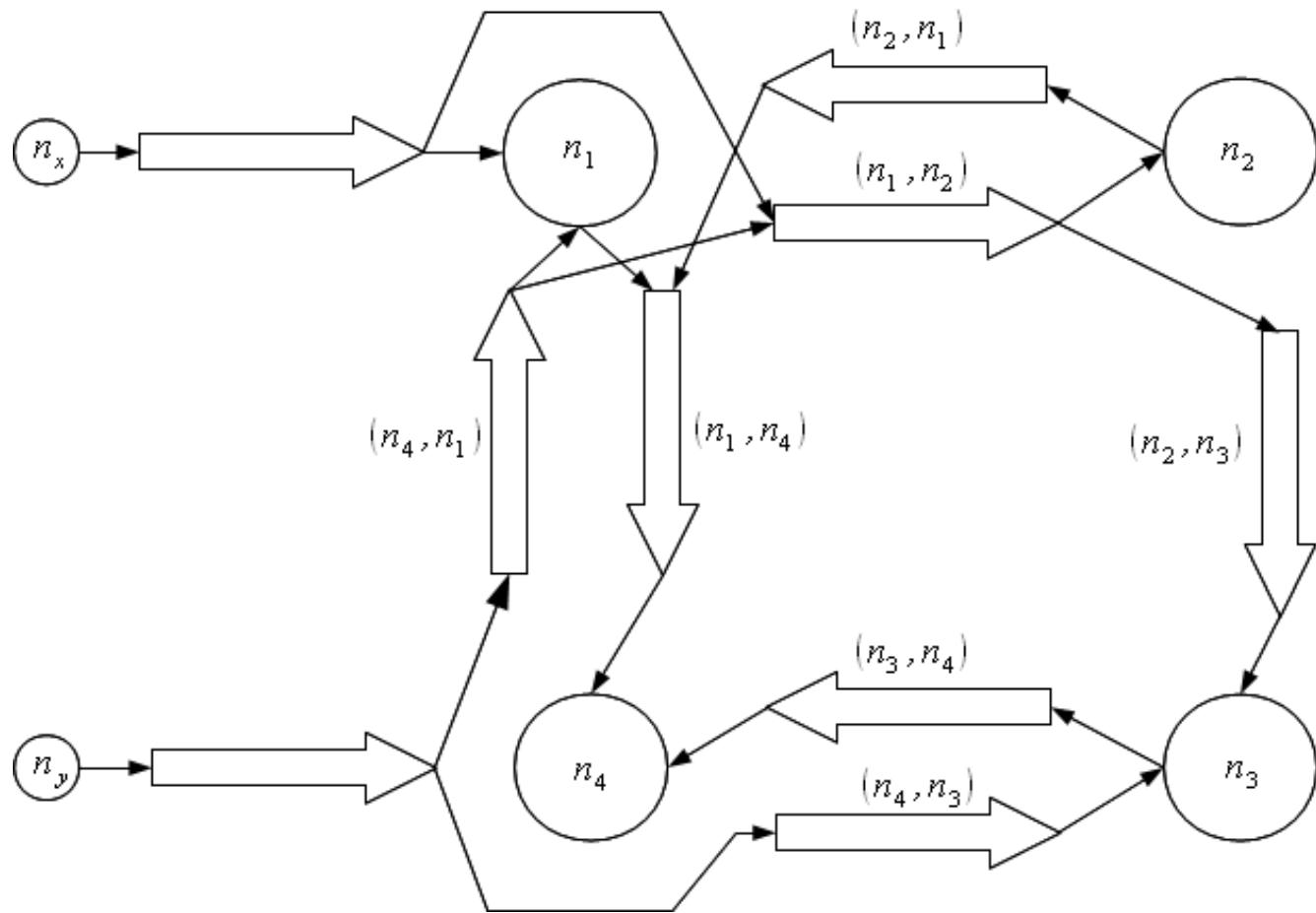
- binary flag determining interconnection between a link (p,n) and link (n,s) :

$$\forall n, p, s \in N \wedge p \in \text{Pred}(n) \wedge s \in \text{Succ}(n) : \exists R_n(p, s) \in \{0,1\}$$

- binary flag determining interconnection between an input node n and link (n,s) :

$$\forall n \in N_i, \forall s \in N \wedge s \in \text{Succ}(n) : \exists S_n(s) \in \{0,1\}$$

FPNN: an example



- For every neural resource there is a set of positive flags:

$$\forall n \in N : \exists \text{set}r_n = \{r_n(p) | r_n(p) = 1 \wedge p \in \text{Pred}(n)\}$$

$$\forall n \in N : \exists \text{sets}_n = \{s_n(s) | s_n(s) = 1 \wedge s \in \text{Succ}(n)\}$$

$$\forall (n, p) \in E : \exists \text{setR}_n(p) = \{R_n(p, s) | R_n(p, s) = 1 \wedge s \in \text{Pred}(n)\}$$

$$\forall n \in N_i : \exists \text{setS}_n(s) = \{S_n(s) | S_n(s) = 1 \wedge s \in N\}$$

- Let N be a set of nodes (**activators**) and E a set of oriented edges (**links**). We say that graph (N,E) is an (standart/light) **grid FPNN** if the following restrictions hold:

- (N,E) is an FPNN
- a limited number of connected preceding links:

$$\forall n \in N : |\text{setr}_n| \geq 2$$

- a limited number of connected successive links:

$$\forall n \in N : |\text{sets}_n| \in (0,1)$$

$$\forall n \in N_i : |\text{setS}_n| \in (0,1)$$

$$\forall (n,p) \in E : |\text{setR}_n(p)| \in (0,1,2)$$

- a limited number of iterations: $\forall n \in N : a_n \geq 2$

- a limited number of inputs: $\forall n \in N_i : c_n = 1$

- for the **standard** grid FPNN:

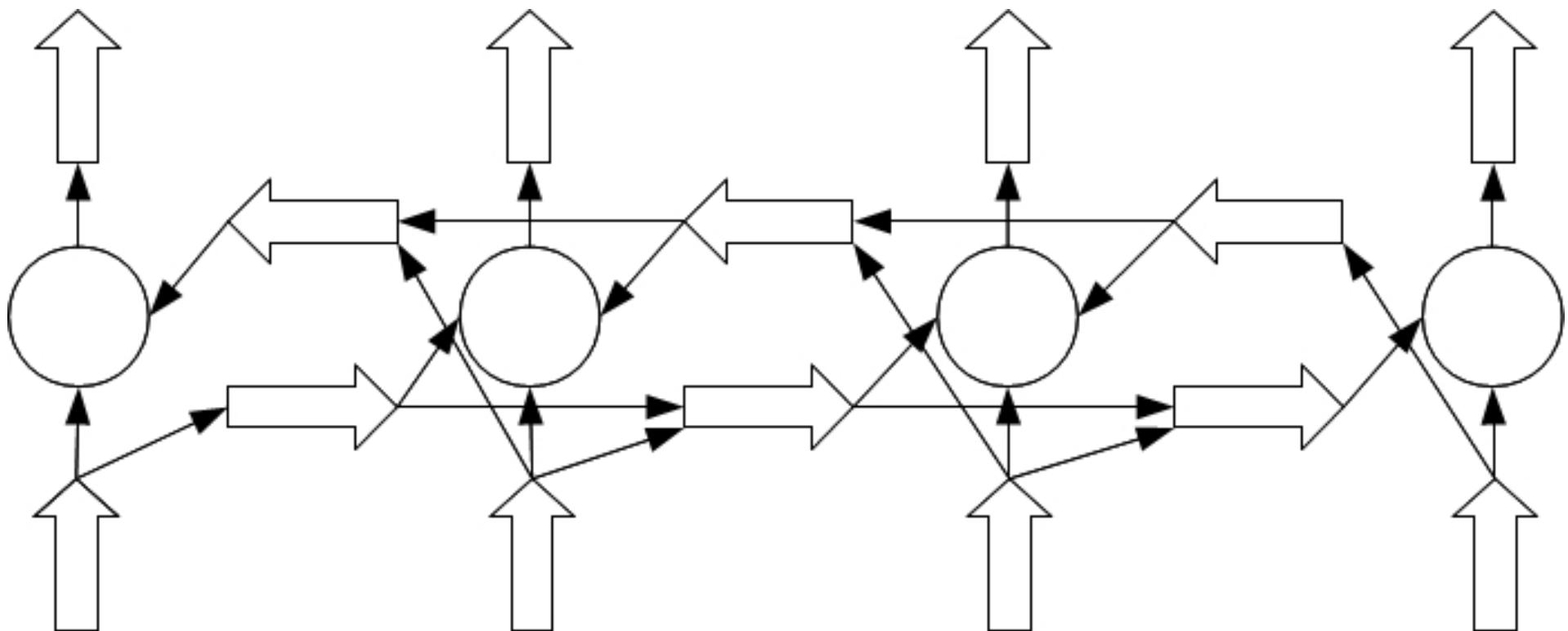
- an affine operator for every predecessor in every link:

$$\forall (r, n) \in E \wedge p \in \text{Pred}(n) : \exists \alpha_{(p, n)} = W_n(p) \cdot x + T_n(p)$$

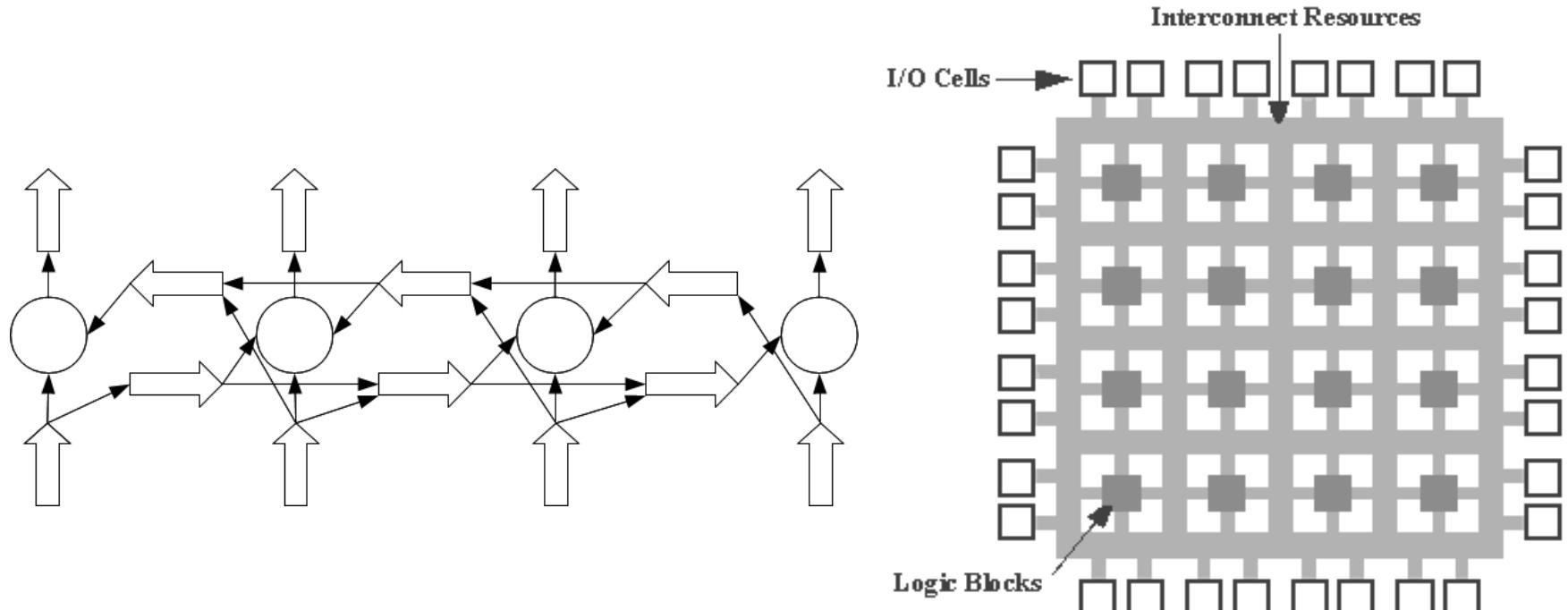
- for the **light** grid FPNN:

- only one affine operator for every link:

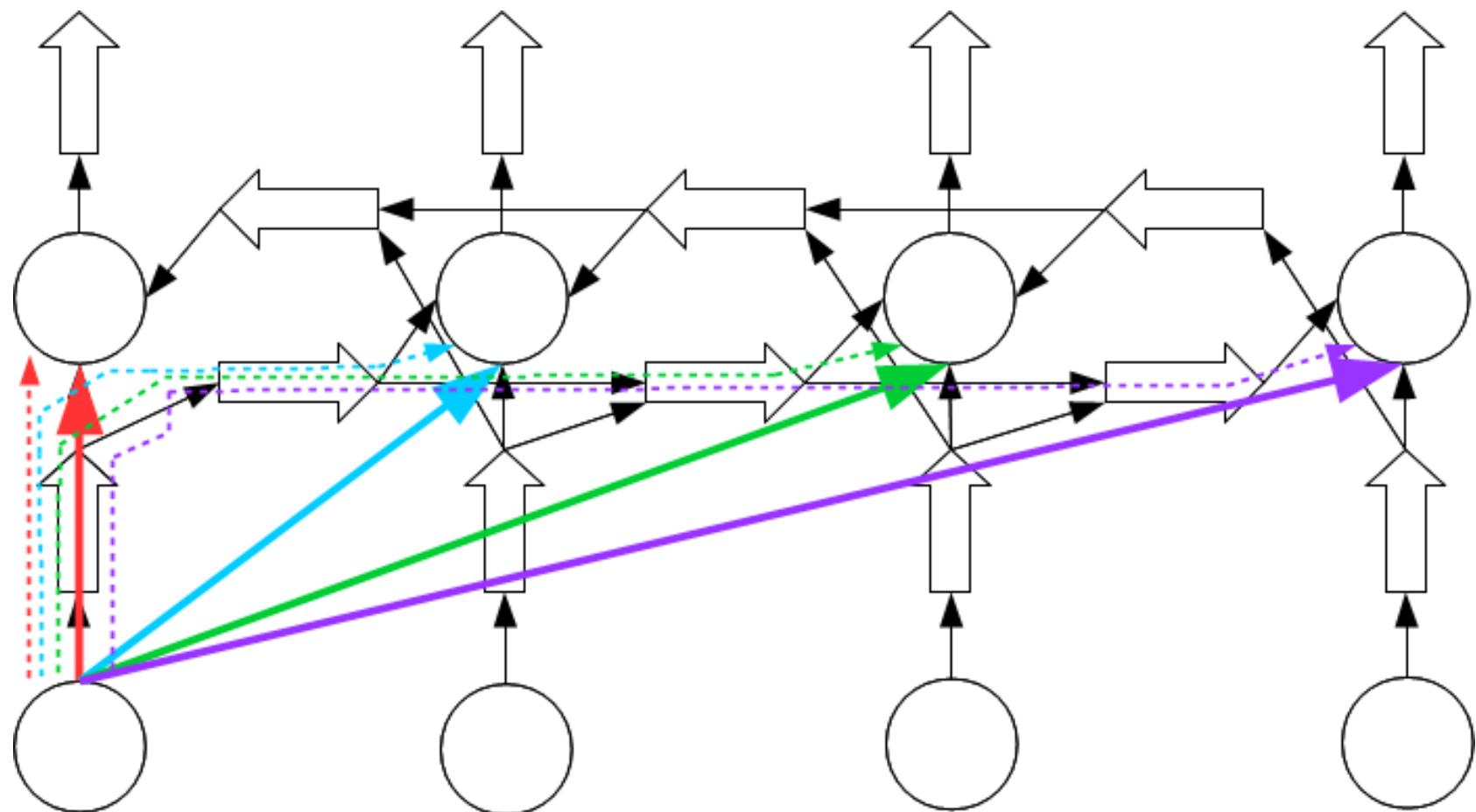
$$\forall p, n \in N \wedge (p, n) \in E : \exists \alpha_{(p, n)} = W \cdot x$$



FPNN: the grid FPNN example



| FPNN: the grid FPNN example



1)Chain of links

- A sequence of links bordered by links having no predecessors (the beginning of the chain) or no successors (the end of the link).

2)Path between links

- Is a subset of chain of links.

0)Declarations

- a) Neural Network is a pair $(\text{Neurons}, \text{Synapses})$
- b) (N, E) is a grid FPNN.
- c) Functions:
 - Mapping neurons to activators: $\text{NeuToAct}: \text{Neurons} \rightarrow N$
 - Inverse mapping: $\text{ActToNeu} = \text{NeuToAct}^{-1}$
 - Source and destination neuron of a synapse:

$\text{srcNeuron}: \text{Synapses} \rightarrow \text{Neurons}$

$\text{dstNeuron}: \text{Synapses} \rightarrow \text{Neurons}$

- Sorting by path length: $\text{sortByLength}: E^n \rightarrow E^n$
- Find path: $\text{findPath}: E \times E \rightarrow E^n, n \geq 1$
- Path's first node : $\text{firstNodeOf}: E^n \rightarrow E$

d) Sets

- Set of synapses ending in the link:
$$\forall u \in E: \exists \text{approximatedSynapses}_u \subset \text{Synapses}$$
- Set of synapses passing through the link:
$$\forall u \in E: \exists \text{passingSynapses}_u \subset \text{Synapses}$$
- Set of connected preceding activators:
$$\forall u \in E: \exists \text{connPred}_u \subset N$$

e) Variables

➤ Synapses weight:

$$\forall syn \in Synapses : \exists syn_{value} \in \mathbb{R}$$

➤ Partial product approximating the synapse:

$$\forall syn \in Synapses : \exists syn_{product} \in \mathbb{R}$$

➤ Last multiplicand of approximation product:

$$\forall syn \in Synapses : \exists syn_{multiplicand} \in \mathbb{R}$$

f) Others

➤ $paths \subset E^*$ is an ordered collection of paths

1) Initialization

- a) Construct a graph of the link connection.
- b) Separate components of the graph and store them in CONN.
- c) Select the links intended to have a non-redundant security and store them in SEL.

2) Mapping of links

$\forall (U, V) \in \text{CONN}:$

a) Expansion of predecessors *approximatedSynapses*:

$\forall (n, s) \in \text{SEL} \wedge (n, s) \in U:$

$\forall (p, n) \in \text{Pred}(n) \wedge R_n(p, s):$

$\forall syn \in \text{approximatedSynapses}_{(n, s)} \wedge \exists act \in \text{connPred}_{(p, n)} \wedge$
 $\wedge \text{ActToNeu}(act) = \text{srcNeuron}(syn) \Rightarrow$
 $\Rightarrow syn \in \text{approximatedSynapses}_{(p, n)}$

b) Paths determining:

$$\text{first} = \{u \in U \mid \deg^+(u) = 0\}$$

$$\text{last} = \{u \in U \mid \deg^-(u) = 0\}$$

$$\text{paths} = \emptyset$$

$\forall u \in \text{first} \wedge \forall v \in \text{last} :$

$$p = \text{findPath}(u, v)$$

$$\text{paths} = \text{paths} \cup 2^p$$

$$\text{sortByLength}(\text{paths})$$

c)mapping path by path

$\forall r \in paths:$

i. Multiplicands computation:

$$(p, n) = firstNodeOf(r)$$

$\forall synapse \in approximatedSynapses_{(p, n)}:$

$$syn_{multiplicand} = \frac{syn_{value}}{syn_{product}}$$

$$\begin{aligned} W_n(p) &= syn_{multiplicand} \Leftrightarrow \\ &\Leftrightarrow srcNeuron(syn) = p \wedge \\ &\wedge dstNeuron(syn) = n \end{aligned}$$

(standard)

$$W_{(p, n)} = \frac{\sum_{syn \in approximatedSynapses_{(p, n)}} syn_{multiplicand}}{|approximatedSynapses_{(p, n)}|}$$

(light)

ii. Updating the products using the new parameters:

$\forall synapse \in passingSynapses_{(p,n)} :$

$$synapse_{product} = synapse_{product} * W_n(p) \Leftrightarrow \\ \Leftrightarrow srcNeuron(syn) = p \wedge dstNeuron(syn) = n \quad (\text{standard})$$

$$synapse_{product} = synapse_{product} * W_{(p,n)} \quad (\text{light})$$

iii. Deleting finished node from the path:

$$r = r \setminus \{(p, n)\}$$

$$paths = \{p | p \in paths \wedge p \neq \emptyset\}$$

1) Declarations

a) Length of rest of a chain: $restOfChain: E^n \times E \rightarrow \mathbb{N}$

b) Set of chains a link belongs to: $chains: E \rightarrow \{E^n\}^m$

c) Fault tolerant ranking: $\forall link \in E: \exists rank_{link} \in \mathbb{R}$

2) Algorithm

$\forall link \in E:$

$$rank 1 = \sum_{syn \in approximatedSynapses_{link}} syn_{value}$$

$$rank 2 = \sum_{chain \in chains(link)} restOfChain(chain, link)$$

$$rank_{link} = rank 1 + rank 2$$

- The FPNN is an approximation of artificial neural networks.
- It is usable for resource saving implementation in the FPGAs.
- Especially grid FPNNs show the resources saving properties.
- Trained artificial neural networks can be directly mapped onto FPNNs.
- The redundancy free fault tolerance can be established during the mapping.
- The links can be ranked with relation to their importance in the circuit.

Bibliography

- Girau, B.: FPNA: Applications and Implementations. In FPGA Implementations of Neural Networks, editation A. R. Omondi; J. C. Rajapakse, Springer US, 2006, ISBN 978-0-387-28487-3, p.81-136, 10.1007/0-387-28487-7-4.
URL <http://dx.doi.org/10.1007/0-387-28487-7-4>
- Martin Krcma: The neural networks acceleration in FPGA, bachelor's thesis, Brno, FIT BUT in Brno, 2012
- Martin Krcma: The neural networks acceleration in FPGA, master's thesis, Brno, FIT BUT in Brno, 2012
- http://en.wikipedia.org/wiki/Field-programmable_gate_array
- <http://home.mit.bme.hu/~szedo/FPGA/fpgahw.htm>

Thank you for your attention