# Sum-Product Networks

Jhonatan Oliveira

Research Tasks A, C, and D

# Outline

# The Paper

---

## Sum-Product Networks: A New Deep Architecture

---

**Hoifung Poon** and **Pedro Domingos**
Computer Science & Engineering
University of Washington
Seattle, WA 98195, USA
{hoifung,pedrod}@cs.washington.edu

### Abstract

The key limiting factor in graphical model inference and learning is the complexity of the partition function. We thus ask the question: what are general conditions under which the partition function is tractable? The answer leads to a new kind of deep architecture, which we call *sum-

the uniform distribution over vectors with an even number of 1's.) Second, inference is still exponential in the worst case. Third, the sample size required for accurate learning is worst-case exponential in scope size. Fourth, because learning requires inference as a subroutine, it can take exponential time even with fixed scopes (unless the partition function is a known constant, which requires restricting the potentials to be conditional probabilities).

# Motivation

- Best Paper Award in UAI'11

- Cited 162 times

- Significant improvement using deep learning

**Sum-Product Networks: A New Deep Architecture**

**Hoifung Poon** and **Pedro Domingos**
Computer Science & Engineering
University of Washington
Seattle, WA 98195, USA
{hoifung,pedrod}@cs.washington.edu

**Abstract**

The key limiting factor in graphical model inference and learning is the complexity of the partition function. We thus ask the question: what are general conditions under which the partition function is tractable? The answer leads to a new kind of deep architecture, which we call *sum-* the uniform distribution over vectors with an even number of 1's.) Second, inference is still exponential in the worst case. Third, the sample size required for accurate learning is worst-case exponential in scope size. Fourth, because learning requires inference as a subroutine, it can take exponential time even with fixed scopes (unless the partition function is a known constant, which requires restricting the potentials to be conditional probabilities).

# Motivation

**Sum-Product Networks: A New Deep Architecture**

**Hoifung Poon** and **Pedro Domingos**
Computer Science & Engineering
University of Washington
Seattle, WA 98195, USA
{hoifung,pedrod}@cs.washington.edu

**Abstract**

The key limiting factor in graphical model inference and learning is the complexity of the partition function. We thus ask the question: what are general conditions under which the partition function is tractable? The answer leads to a new kind of deep architecture, which we call *sum-* the uniform distribution over vectors with an even number of 1's.) Second, inference is still exponential in the worst case. Third, the sample size required for accurate learning is worst-case exponential in scope size. Fourth, because learning requires inference as a subroutine, it can take exponential time even with fixed scopes (unless the partition function is a known constant, which requires restricting the potentials to be conditional probabilities).

**Shallow vs. Deep Sum-Product Networks**

**Olivier Delalleau**
Department of Computer Science and Operation Research
Université de Montréal
delallea@iro.umontreal.ca

**Yoshua Bengio**
Department of Computer Science and Operation Research
Université de Montréal
yoshua.bengio@umontreal.ca

**Abstract**

We investigate the representational power of sum-product networks (computation networks analogous to neural networks, but whose individual units compute either products or weighted sums), through a theoretical analysis that compares deep (multiple hidden layers) vs. shallow (one hidden layer) architectures. We prove there exist families of functions that can be represented much more efficiently with a deep network than with a shallow one, i.e. with substantially fewer hidden units. Such results were not available until now, and contribute to motivate recent research involving learning of deep sum-product networks, and more generally motivate research in Deep Learning.

## Introduction and prior work

any learning algorithms are based on searching a family of functions so as to identify one memb said family which minimizes a training criterion. The choice of this family of functions and ho mbers of that family are parameterized can be a crucial one. Although there is no universall timal choice of parameterization or family of functions (or "architecture"), as demonstrated b no-free-lunch results [37], it may be the case that some architectures are appropriate (or inap

# Motivation

## Resource Constrains

- Deep learning:

  - processing

  - data management

- Deep architecture

  - storage

  - inference



nVIDIA.

GPU ACCELERATED DEEP LEARNING
WITH CUDNN

# Background

## Join Probability Distribution

- A multivariate function over a finite set of variables

- Assigns a real number between 0 and 1 to each configuration (combination of variable's values) of the variables

- Summing all assigned real numbers yields 1

| X1 | X2 | P(X1,X2) |
|----|----|----------|
| 0 | 0 | 0.1 |
| 0 | 1 | 0.3 |
| 1 | 0 | 0.5 |
| 1 | 1 | 0.1 |

# Background

| A | B | P(X1,X2) | |
|---|---|---|---|
| 0 | 0 | 0.1 | $\lambda_{A_0} \lambda_{B_0}$ |
| 0 | 1 | 0.3 | $\lambda_{A_0} \lambda_{B_1}$ |
| 1 | 0 | 0.5 | $\lambda_{A_1} \lambda_{B_0}$ |
| 1 | 1 | 0.1 | $\lambda_{A_1} \lambda_{B_1}$ |

## Network Polynomial

$$f = 0.1\lambda_{A_0}\lambda_{B_0} + 0.3\lambda_{A_0}\lambda_{B_1} + 0.5\lambda_{A_1}\lambda_{B_0} + 0.1\lambda_{A_1}\lambda_{B_1}$$

# Background

| A | B | P(X1,X2) | |
|---|---|---|---|
| 0 | 0 | 0.1 | $\lambda_{A_0} \lambda_{B_0}$ |
| 0 | 1 | 0.3 | $\lambda_{A_0} \lambda_{B_1}$ |
| ~~1~~ | ~~0~~ | ~~0.5~~ | $\lambda_{A_1} \lambda_{B_0}$ |
| ~~1~~ | ~~1~~ | ~~0.1~~ | $\lambda_{A_1} \lambda_{B_1}$ |

**Network Polynomial**

$$f \;=\; 0.1\lambda_{A_0}\lambda_{B_0} + 0.3\lambda_{A_0}\lambda_{B_1} + 0.5\lambda_{A_1}\lambda_{B_0} + 0.1\lambda_{A_1}\lambda_{B_1}$$

# Background



**Network Polynomial**

$$f = 0.1\lambda_{A_0}\lambda_{B_0} + 0.3\lambda_{A_0}\lambda_{B_1} + 0.5\lambda_{A_1}\lambda_{B_0} + 0.1\lambda_{A_1}\lambda_{B_1}$$

# Background



**Sum-Product Network**
a directed acyclic graph formed by
sum and product nodes.

# Critic

## Good Points

- Novel model

- Well written

- Logical flow

- Experiments

- Relate to other field

### 6 SUM-PRODUCT NETWORKS AND THE CORTEX

The cortex is composed of two main types of cells: pyramidal neurons and stellate neurons. Pyramidal neurons excite the neurons they connect to; most stellate neurons inhibit them. There is an interesting analogy between these two types of neurons and the nodes in SPNs, particularly when MAP inference is used. In this case the network is composed of max nodes and sum nodes (logs of products). (Cf. Riesenhuber and Poggio [23], which also uses max and sum nodes, but is not a probabilistic model.) Max nodes are analogous to inhibitory neurons in that they select the highest input for further propagation. Sum nodes are analogous to excitatory neurons in that they compute a sum of their inputs. In SPNs the weights are at the inputs of max nodes, while the analogy with the cortex suggests having them at the inputs of sum (log product) nodes. One can be mapped to the other if we let max nodes ignore their children's weights and consider only their values. Possible justifications for this include: (a) it potentially reduces computational cost by allowing max nodes to be merged; (b) ignoring priors may improve discriminative performance [11]; (c) priors may be approximately encoded by the number of units representing the same pattern, and this may facilitate online hard EM learning. Unlike SPNs, the cortex has no single root node, but it is straightforward to extend SPNs to have multiple roots, corresponding to simultaneously computing multiple distributions with shared structure. Of course, SPNs are still biologically unrealistic in

# Critic

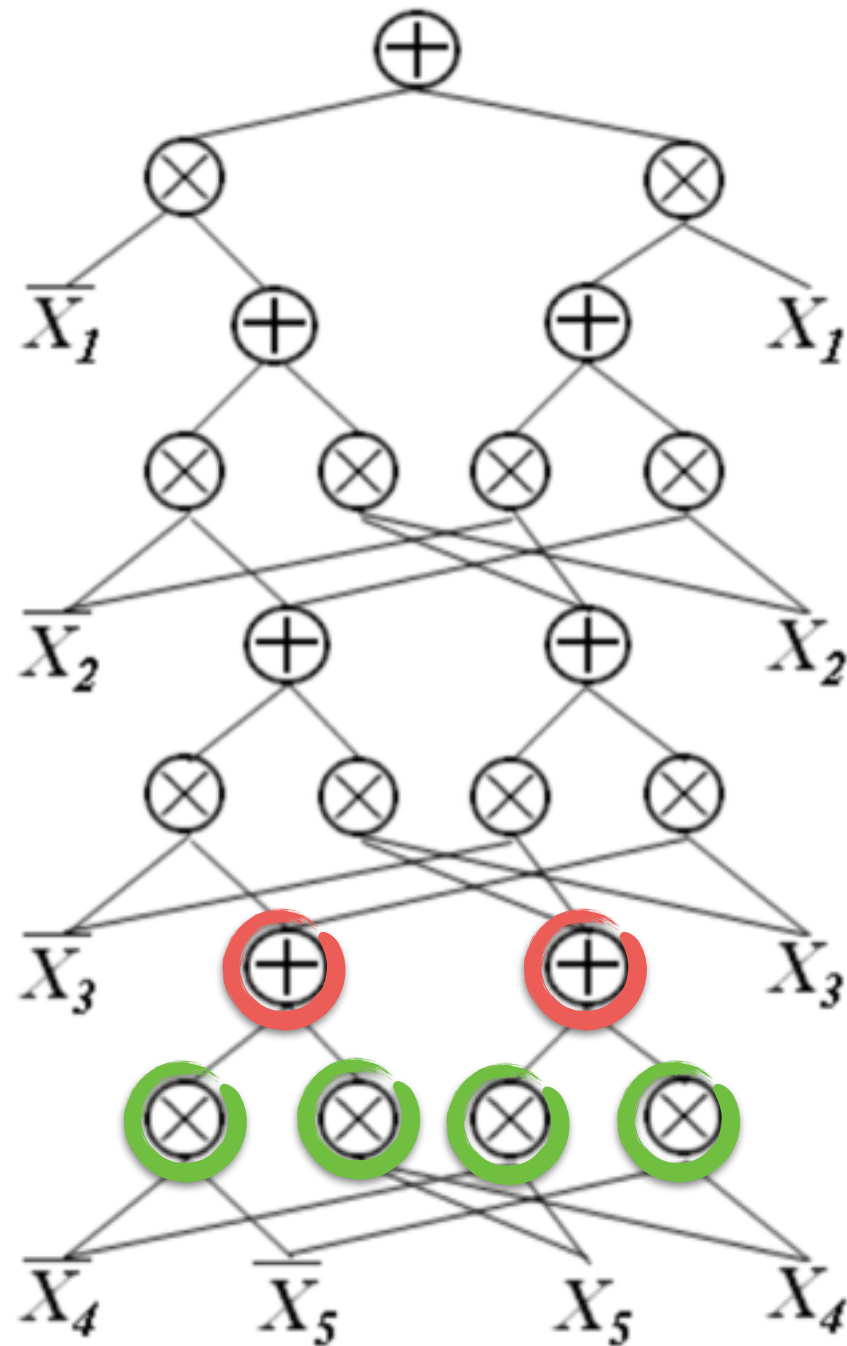**Points to improve**
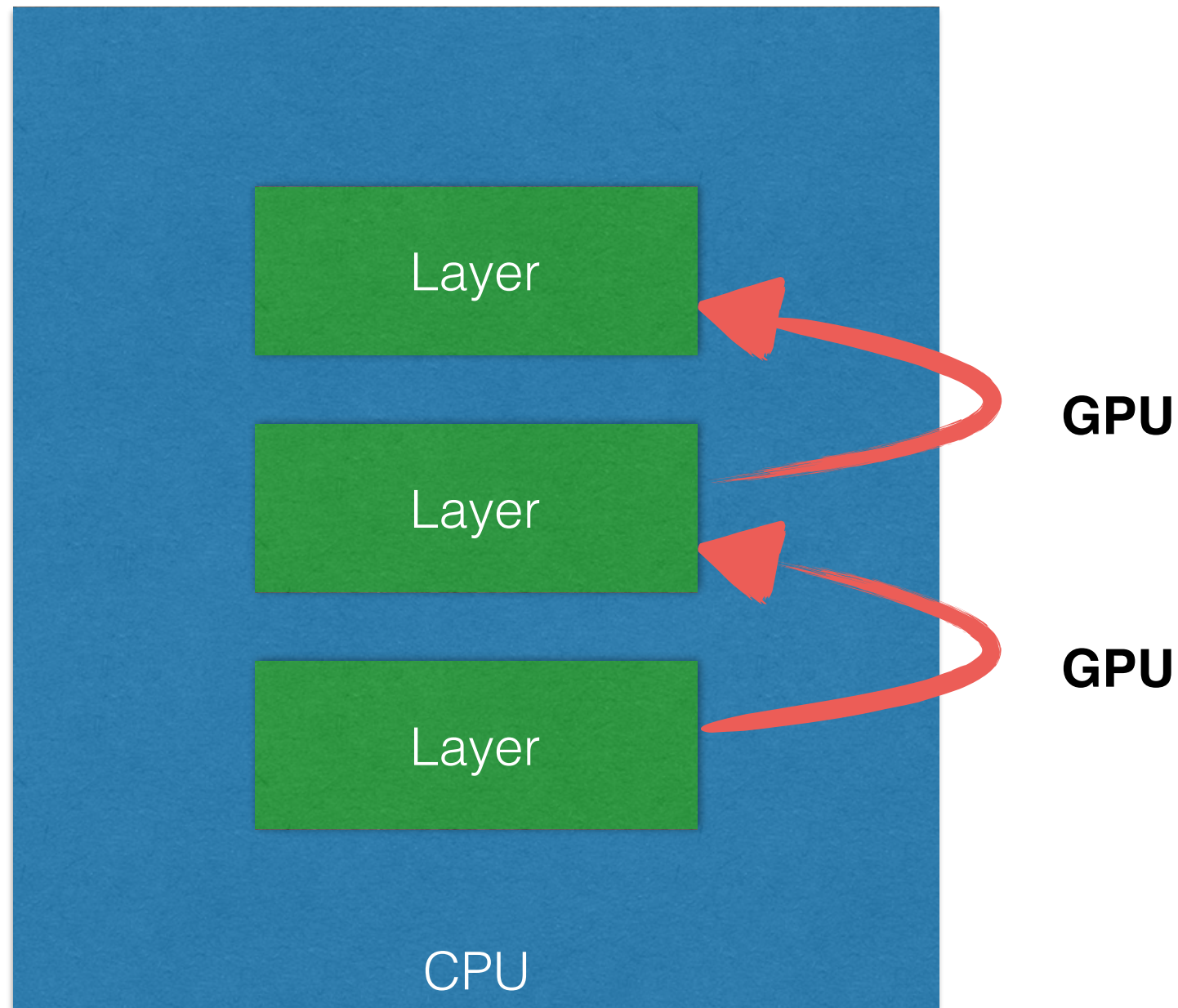
- Related works

- Background

# The Problem

# Parallel Computation
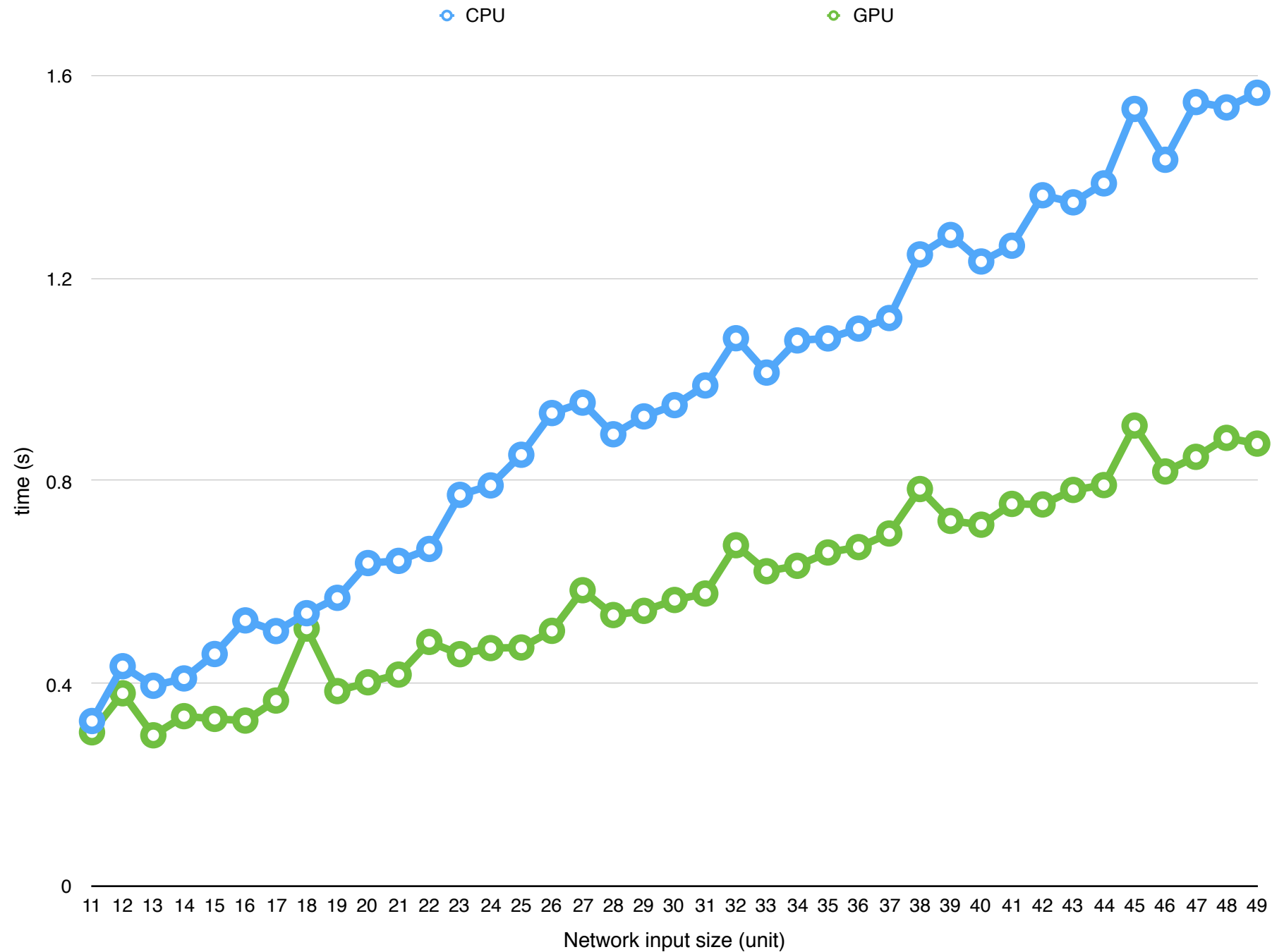
# Possible Issues

# The Comparison

# Implementation

# Implementation

```python
import numpy as np
from keras.models import Sequential
from keras.layers.core import Dense, Dropout, Layer, Activation
import time
import tensorflow as tf

f = open("results.csv", "w")

for i in range(10, 50):

    print("--> Running at %d <--" % i)
    INPUT_SIZE = i
    OUTPUT_SIZE = INPUT_SIZE
    nb_class = 3

    # global parameters
    batch_size = 128
    nb_epoch = 40

    np.random.seed(123)

    X_train = np.random.rand(INPUT_SIZE, nb_class)
    Y_train = np.random.rand(OUTPUT_SIZE, nb_class)

    X_test = np.random.rand(INPUT_SIZE)
    Y_test = np.random.rand(OUTPUT_SIZE)

    print("--> Building model...")
    model = Sequential()
    model.add(Dense(INPUT_SIZE, input_shape=(nb_class,)))
    model.add(Activation('linear'))
    model.add(Dense(OUTPUT_SIZE))
    model.add(Activation('linear'))

    print("--> Compiling model...")

    start_time = time.time()

    model.compile(loss='categorical_crossentropy', optimizer='rmsprop')

    final_time = time.time()
    diff_time = final_time - start_time
    print("--> Exec time:")
    print(diff_time)

    f.write(str(i)+","+str(diff_time)+","+"\n")

f.close()
```

# CPU x GPU

# Conclusions

- Sum-product networks are a new deep architecture for modelling and inference

- Paper is well done with good ideas and nice presentation

- Inference and learning can be parallelized per layer

- GPUs can speed up SPN learning and inference