

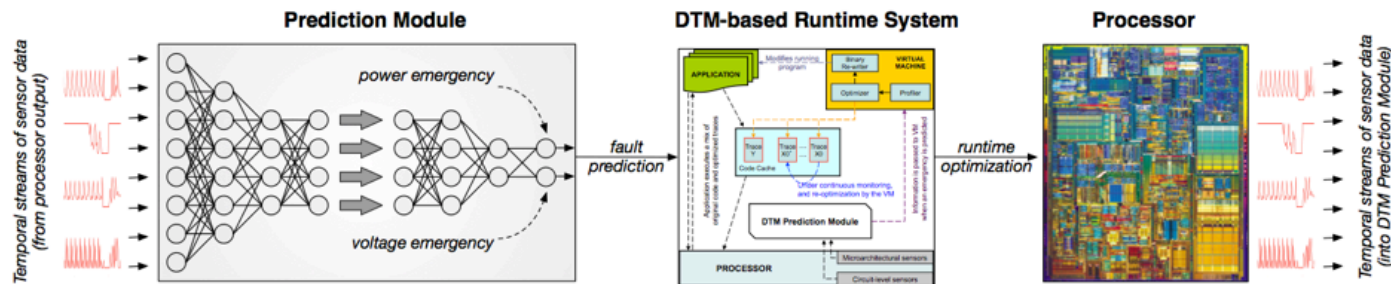
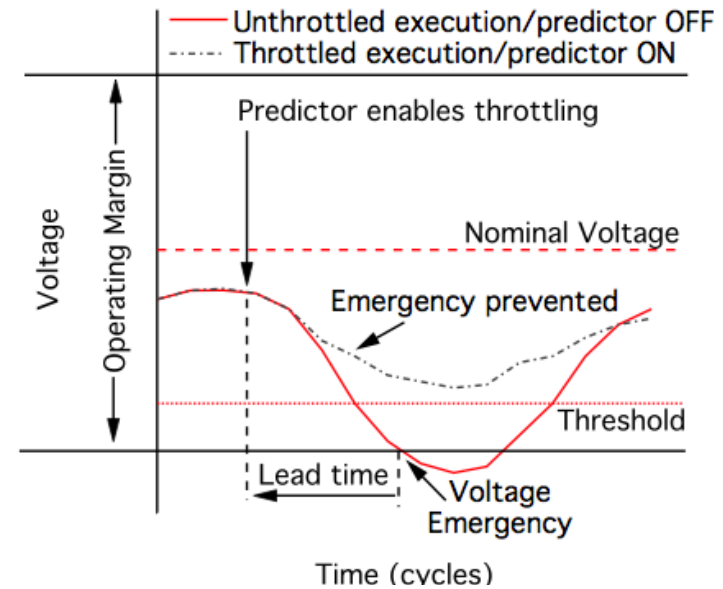
# Event Prediction in Processors using Deep Temporal Models

Tharindu Mathew, Aswin Raghavan, Sek Chai  
SRI International

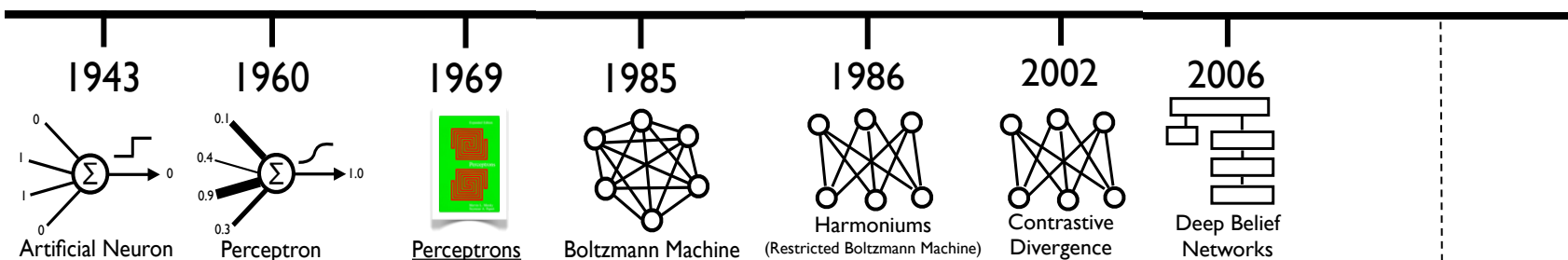
# NSF Resilient Computing

The overall objective of our work is to lay the foundation for a **proactive** computing platform that manages its own fault resiliency and reliability.

The ultimate goal is to eliminate the penalties to system design and operation that arise in the use of incremental circuit techniques and micro-architectural changes that increase system complexity, affect power consumption and cost.

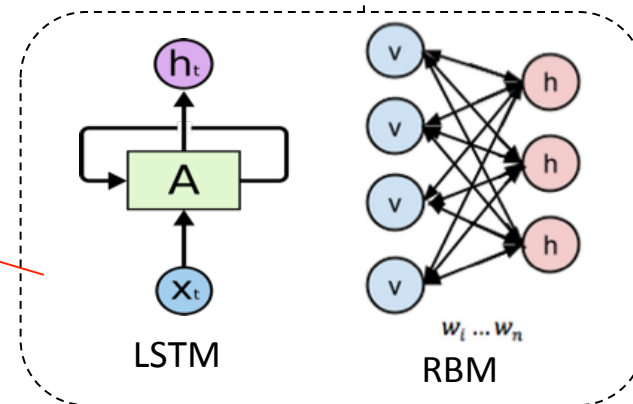


# Brief Taxonomy



## Deep Temporal Models (DTM)

- Are hierarchical and multi-layered
- Address spatio-temporal data
- Typically, recurrent connections.
- Examples: LSTM, RBM



## Deep Temporal Models

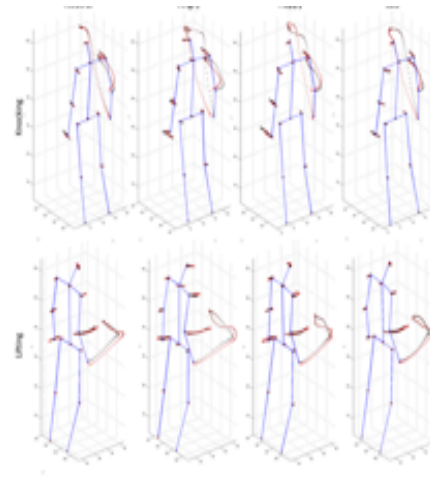
Timeline info loosely based on Yoshua Bengio's talk:

[https://homes.cs.washington.edu/~rcg/talks/BengioDeepArchitectures\\_GensCommentary.pdf](https://homes.cs.washington.edu/~rcg/talks/BengioDeepArchitectures_GensCommentary.pdf)

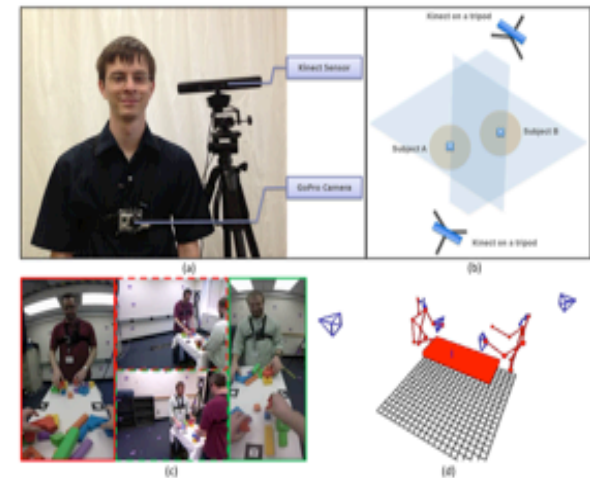
## Applications using DTM (Benchmarks 1D – 4D dimensional data)



Facial Attributes Classification  
(Images 2D/Multi Task)



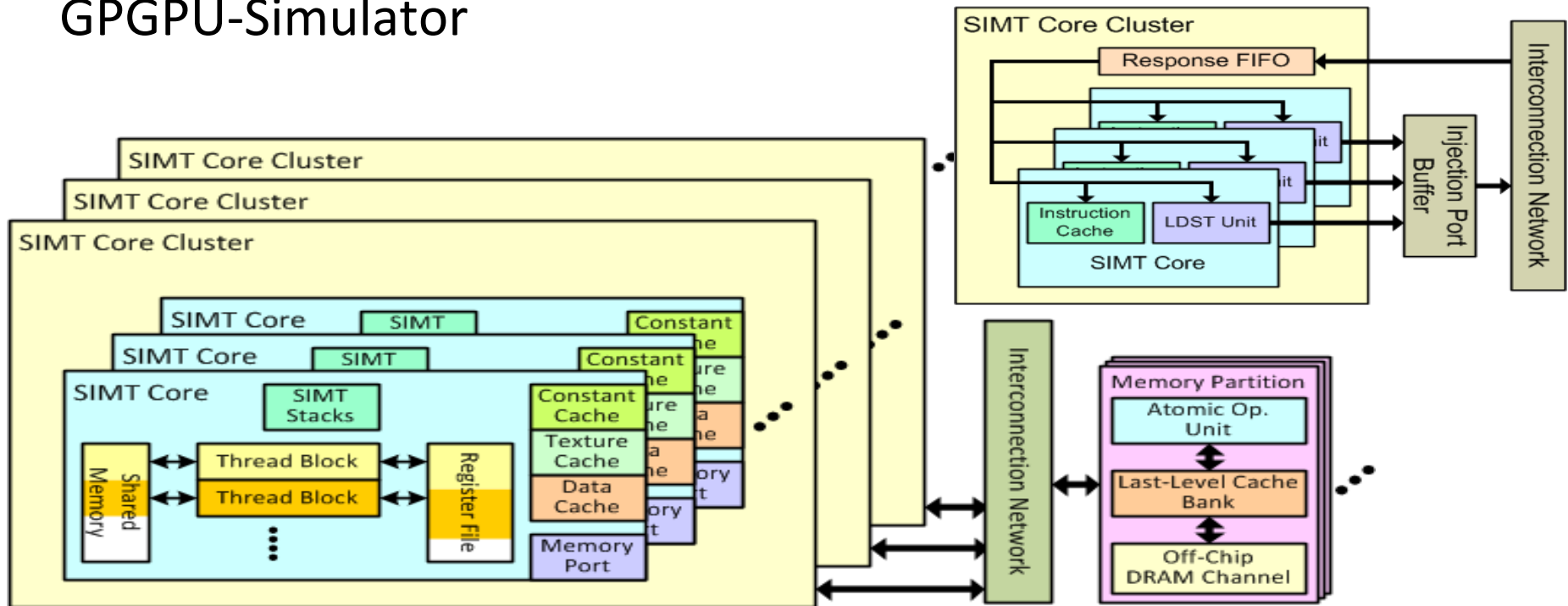
Body Affect Recognition  
(Mocap 3D+t/Multi Task)



Tower Game for Engagement Classification  
(Mocap, Audio, Depth, RGB: N-D+t/Multi Task)

Amer, M. R., Shields, T., Siddique, B., Tamrakar, A., Divakaran, A., & Chai, S. (2018). Deep multimodal fusion: A hybrid approach. International Journal of Computer Vision, 126(2-4), 440-456.

# GPGPU-Simulator



<http://www.gpgpu-sim.org/>

- Kernels is a block of code such as a C function.
- Each core contains many thread blocks that execute SIMT.

## GPGPU-Sim Kernels

	Application	Dwarf	Domain	Problem Sizes
1	<b>Back Propagation (BP)</b>	Unstructured Grid	Pattern Recognition	65536 input nodes
2	<b>Breadth First Search (BF)</b>	Graph Traversal	Graph Algorithms	1000000 nodes
3	<b>CFD Solver (C)</b>	Unstructured Grid	Fluid Dynamics	97k elements
4	<b>Heart Wall Tracking (HW)</b>	Structured Grid	Medical Imaging	609 x 590 pixels/frame
5	<b>HotSpot (HP)</b>	Structured Grid	Physics Simulation	500 x 500 data points
6	<b>Kmeans (K)</b>	Dense Linear Algebra	Data Mining	204800 data points, 34 features
7	<b>Leukocyte Tracking (LK)</b>	Structured Grid	Medical Imaging	219 x 640 pixels/frame
8	<b>LU Decomposition (LU)</b>	Dense Linear Algebra	Linear Algebra	256 x 256 data points
9	<b>MUMmer (M)</b>	Graph Traversal	Bioinformatics	50000 25-character queries
10	<b>Needleman-Wensch (N)</b>	Dynamic Programming	Bioinformatics	2048 x 2048 data points
11	<b>SRAD (SR)</b>	Structured Grid	Image Processing	512 x 512 data points
12	<b>Stream Cluster (SC)</b>	Dense Linear Algebra	Data Mining	65536 points, 256 dimensions

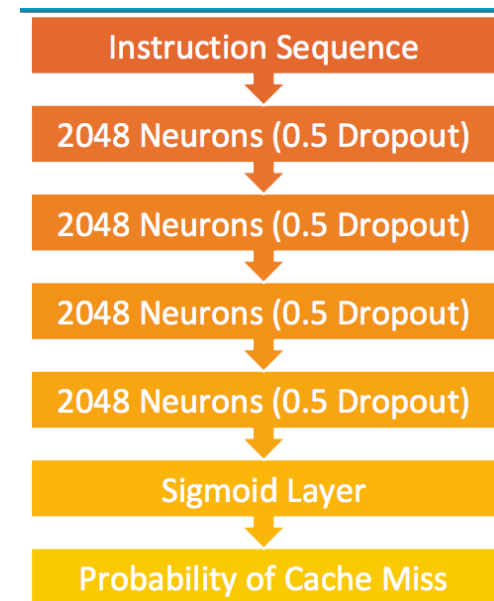
Twelve (12) programs of the Rodinia Benchmark is run through the GPGPU simulator simulating a NVidia GTX 480. Executed instructions across cores and data cache read miss (DC\_RM) data are extracted from all of the programs.

# Time-series analysis using Deep Temporal Models

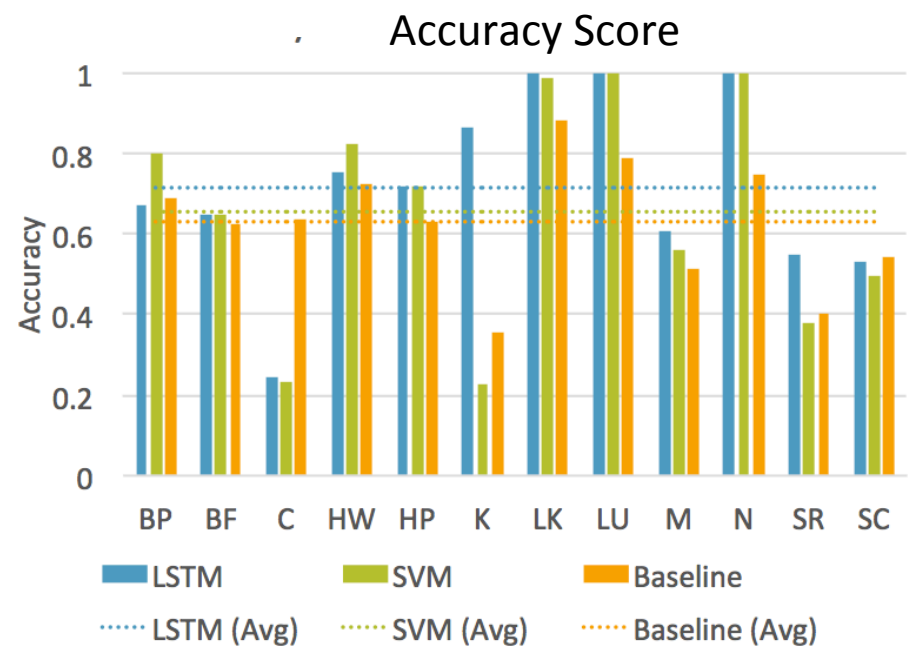
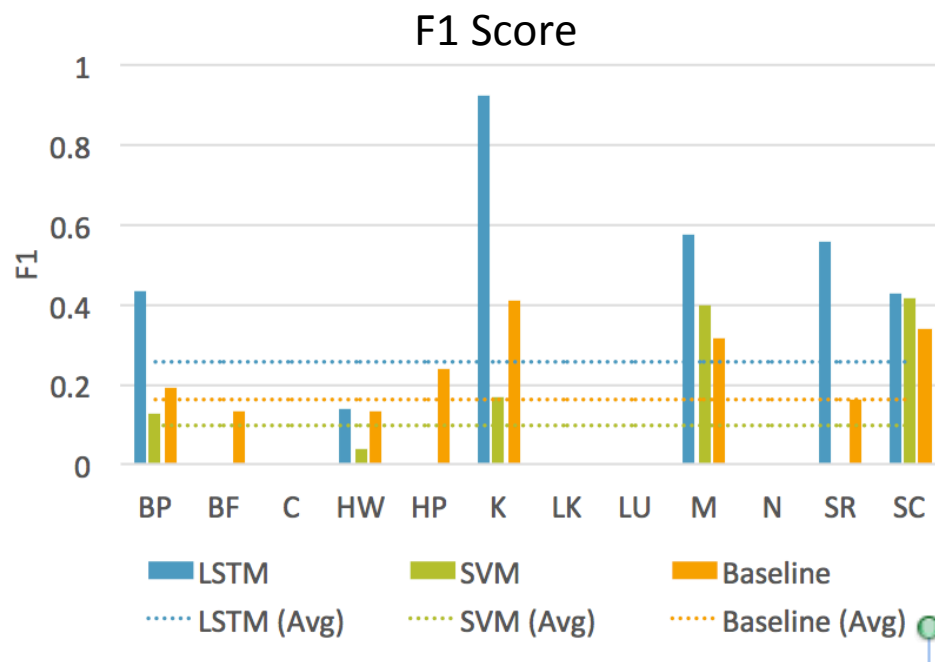
DTM is used as a **sequence predictor**.

The GPU instruction sequences (e.g. add, mov, call) are  $m$  dimensional inputs with the occurrence of a cache miss as the binary output at each execution cycle. The deep temporal model is then trained to look at sequence of instructions, and predict the occurrence of a cache miss,  $n$  cycles ahead.

- For our Deep Temporal model, we use 2048-neuron LSTM cell in 4 successive layers.
- To improve regularization, a drop out ratio of 0.5 is applied after each layer.
- A sigmoid layer outputs the probability of a cache miss.
- A threshold of 0.5 is applied to create a binary output that represents the cache miss.



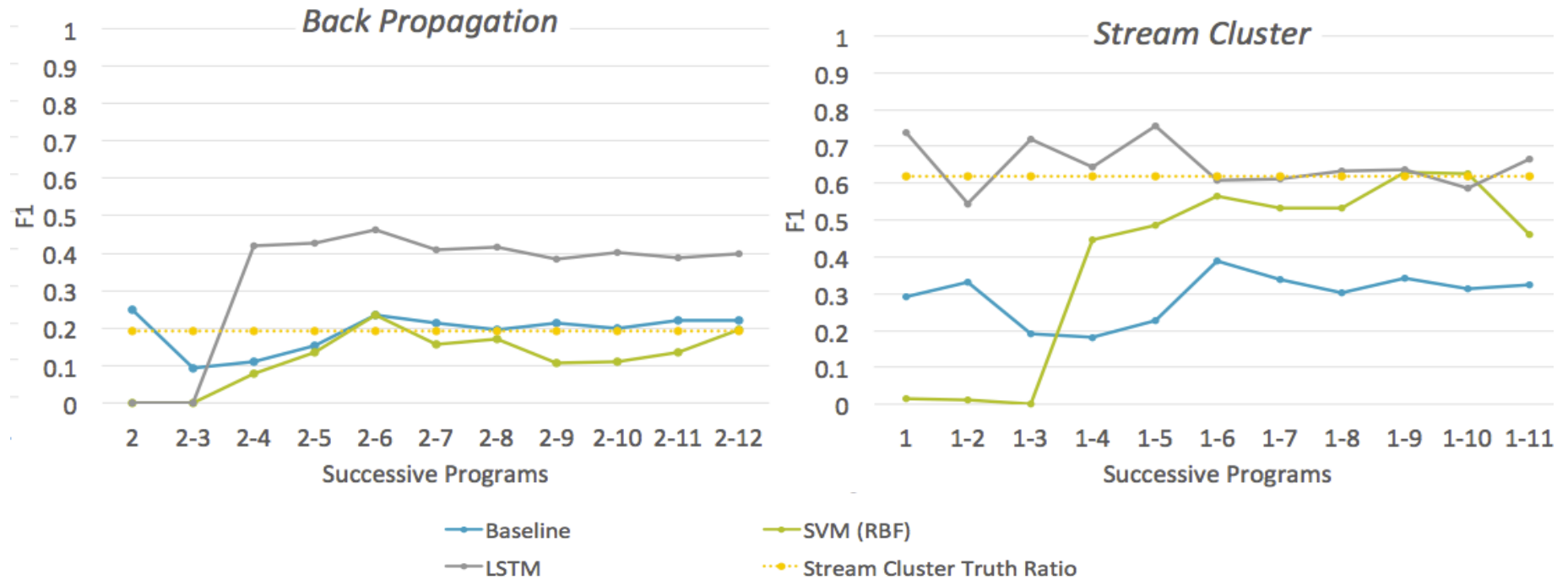
# Predicting Cache Events



Five (5) randomly picked training programs are used as training programs to predict misses of an unseen test program



## Successive cross-dataset prediction



Multiple program instructions are stacked successively for training. e.g. 1-6:  
instructions from BP+BF+C+HW+HP+K are stacked together