

# A FULLY PIPELINED KERNEL NORMALISED LEAST MEAN SQUARES PROCESSOR FOR ACCELERATED PARAMETER OPTIMISATION

---

**Nicholas J. Fraser**, Duncan J.M. Moss, JunKyu Lee, Stephen Tridgell, Craig T. Jin and Philip H.W. Leong

September 3, 2015

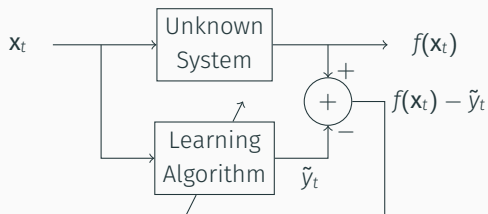
CELab, University of Sydney, Australia

## Contributions:

- A fully pipelined core implementing KNLMS.
- A complete PCIe system.
- Large speedups over previous designs.
- Floating point and *fused* arithmetic designs - details in the paper!

# INTRODUCTION: MOTIVATION & AIMS

Machine learning - creating models to fit data.



Two common problems arise:

- How do we decide which algorithm to use?
- Given an algorithm, how do we configure it?

# MOTIVATION: KNLMS CONFIGURATION

Parameter selection can be the difference between a good and bad models.

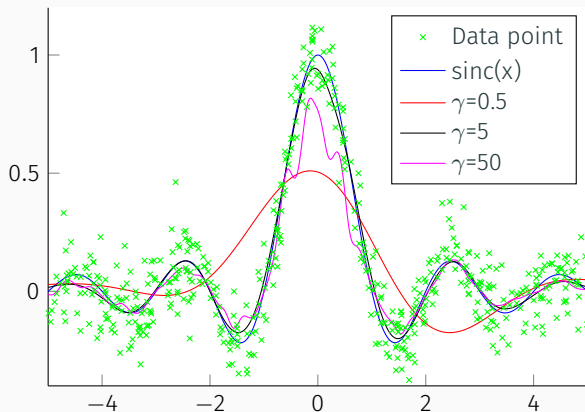


Figure 1: The accuracy of KNLMS while varying a single parameter.

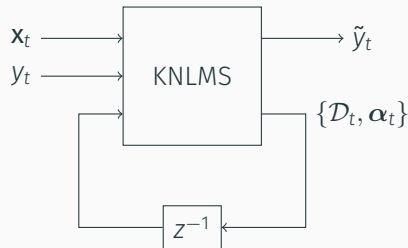
# BACKGROUND: KERNEL METHODS

KNLMS model:

- A dictionary,  $\mathcal{D}$ .
- A vector of weights,  $\alpha$ .
- A kernel function,  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$ .

Prediction calculation (given  $\mathbf{x}_t$ ):

$$\tilde{y}_t = \sum_{i=0}^N \alpha_i \kappa(\mathbf{x}_t, \tilde{\mathbf{x}}_i) = \mathbf{k}^T \boldsymbol{\alpha} .$$



The expressions are recursive - this creates a dependency problem!

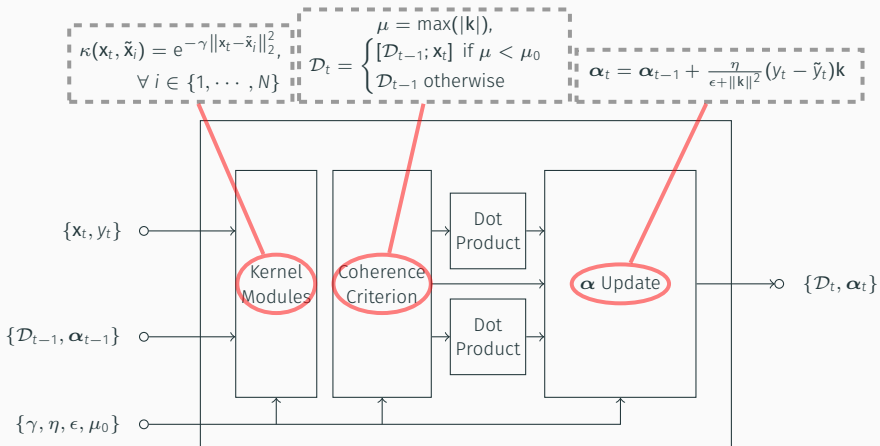
Avoid the dependency problem by reordering the loop.

```
for (parameters) {  
  for (examples) {  
    learn_model();  
  }  
}
```

→  
→

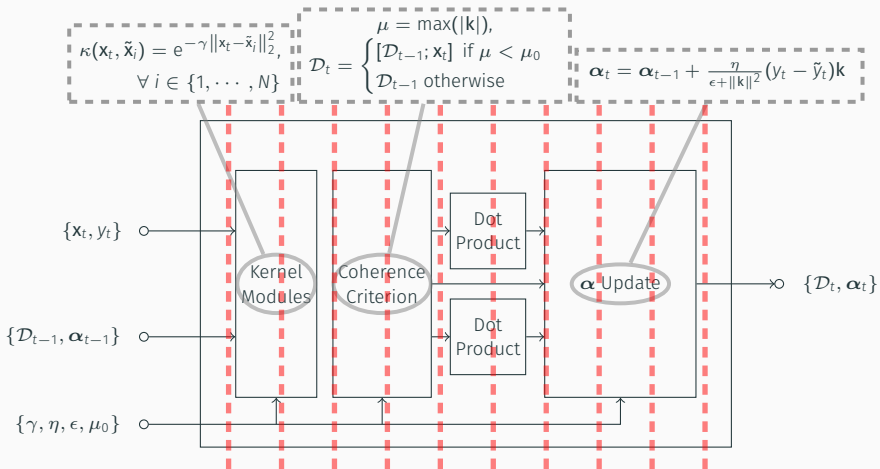
```
for (examples) {  
  for (parameters) {  
    learn_model();  
  }  
}
```

# ARCHITECTURE: THE KNLMS CORE



The core calculates a **non-recursive** portion of KNLMS.

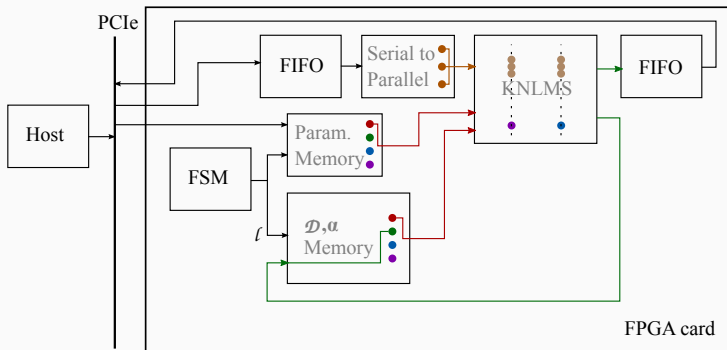
# ARCHITECTURE: PIPELINING THE CORE



~210 pipeline stages were needed to achieve ~300MHz.



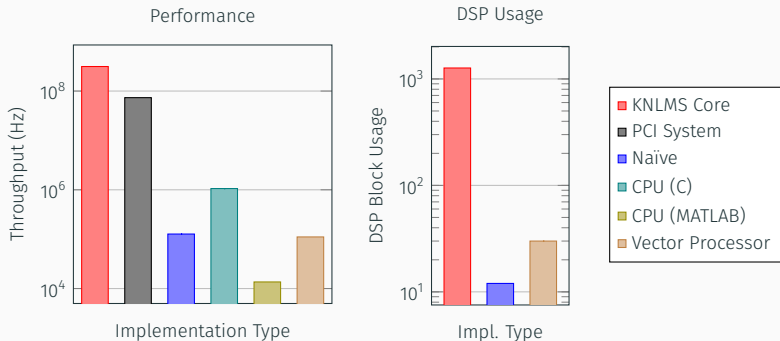
# ARCHITECTURE: PCI-BASED SYSTEM



On each clock cycle, a different set of parameters is passed in.

# RESULTS: PERFORMANCE

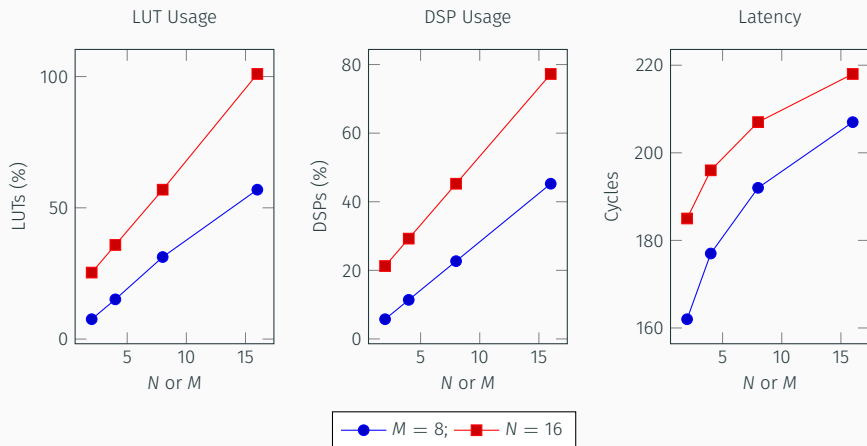
The core achieves speedups of  $300\times/2,800\times$  over a CPU (C) and a vector processor respectively.



The core is capable of learning  $\sim 210$  models at  $\sim 160$ GFLOPS.

# RESULTS: SCALABILITY

- VC707 (Virtex 7):  $N=16, M=8$ . ( $N$  = dictionary entries,  $M$  = feature length)
- Virtex Ultrascale+:  $N=64, M=8$  (estimated).



## CONCLUSION: SUMMARY

- Demonstration of a fully pipelined machine learning core:
  - The core achieves  $\sim 300\times / \sim 2,800\times$  speedups over a CPU and a previous design.
  - A 210 stage pipeline core achieves 160GFLOPS.
  - PCI system achieves  $\sim 70\times / \sim 660\times$  speedups over a CPU and a previous design.
- Floating point and fused arithmetic investigated - details in the paper!
- We hope this design enables embedded and real-time applications for machine learning.

Future work includes:

- Further investigation of precision tradeoffs.
- Reducing the design latency.
- Implementing other machine learning algorithms.
- Using other platforms, such as GPUs and heterogeneous architectures.

THANK YOU. ANY QUESTIONS?

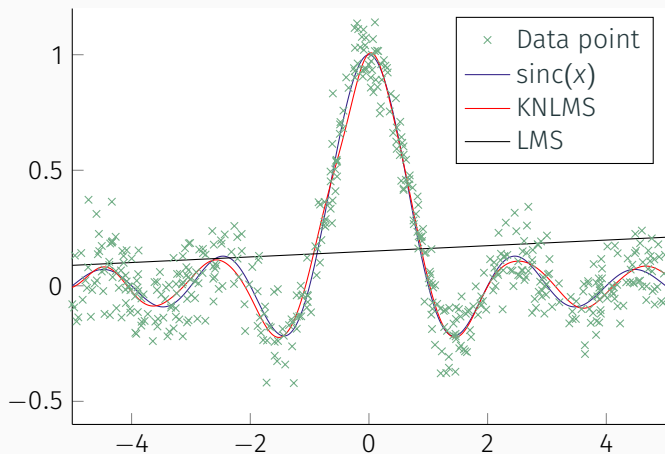
- Y. Engel, S. Mannor, and R. Meir. The kernel recursive least-squares algorithm. *Signal Processing, IEEE Transactions on*, 52(8):2275–2285, 2004.
- M. Jacobsen, Y. Freund, and R. Kastner. RIFFA: A reusable integration framework for FPGA accelerators. In *Field-Programmable Custom Computing Machines (FCCM), 2012 IEEE 20th Annual International Symposium on*, pages 216–219. IEEE, 2012. ISBN 978-0-7695-4699-5.
- Y. Pang, S. Wang, Y. Peng, N. J. Fraser, and P. H. Leong. A low latency kernel recursive least squares processor using FPGA technology. In *FPT*, pages 144–151, 2013.
- J. Platt et al. Sequential minimal optimization: A fast algorithm for training support vector machines. 1998.

- C. Richard, J. C. M. Bermudez, and P. Honeine. Online prediction of time series data with kernels. *Signal Processing, IEEE Transactions on*, 57(3):1058–1067, 2009.
- S. Van Vaerenbergh. Kernel methods toolbox KAFBOX: a Matlab benchmarking toolbox for kernel adaptive filtering. Grupo de Tratamiento Avanzado de Señal, Departamento de Ingeniería de Comunicaciones, Universidad de Cantabria, Spain, 2012.
- B. Widrow and M. J. Hoff. Adaptive switching circuits. In *IRE WESCON Convention Record*, pages 96–104, 1960.



# APPENDICES

# KNLMS Vs LMS



**Figure 2:** Comparison of the ability of LMS and KNLMS the learn  $f(x) = \text{sinc}(x)$ .

**Table 1:** Computational complexity of some different machine learning methods when used in an online setting. Note that usually  $m \leq n \ll N$ .

Algorithm Type	Computational Cost	Complexity
LMS [Widrow and Hoff, 1960]	$\mathcal{O}(m)$	(Simple)
KNLMS [Richard et al., 2009]	$\mathcal{O}(nm)$	(Modest)
KRLS [Engel et al., 2004]	$\mathcal{O}(nm + n^2)$	(Moderate)
SVM [Platt et al., 1998]	$\sim (N) \rightarrow (N^{2.2})$	(High)

Some kernel functions include:

- The polynomial kernel:  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + c)^d$ .
- The Gaussian kernel:  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|_2^2}$ .
- For the Gaussian kernel: a measure of “similarity” between input vectors. The Gaussian kernel was used in this work.

Given a new input/output pair,  $\{\mathbf{x}_t, y_t\}$ , a model update is calculated as follows:

1. Evaluate  $\kappa$  between  $\mathbf{x}_t$  and each entry of  $\mathcal{D}_{t-1}$ , creating a *kernel vector*,  $\mathbf{k}$ .
2. If  $\max(|\mathbf{k}|) < \mu_0$ , add  $\mathbf{x}_t$  to the dictionary, producing  $\mathcal{D}_t$ .
3. Update the weights using:  $\boldsymbol{\alpha}_t = \boldsymbol{\alpha}_{t-1} + \frac{\eta}{\epsilon + \mathbf{k}^T \mathbf{k}} (y_t - \mathbf{k}^T \boldsymbol{\alpha}_{t-1}) \mathbf{k}$ .

How can we choose  $\kappa$ ,  $\mu_0$ ,  $\eta$  and  $\epsilon$ ? We must do a *parameter search*.

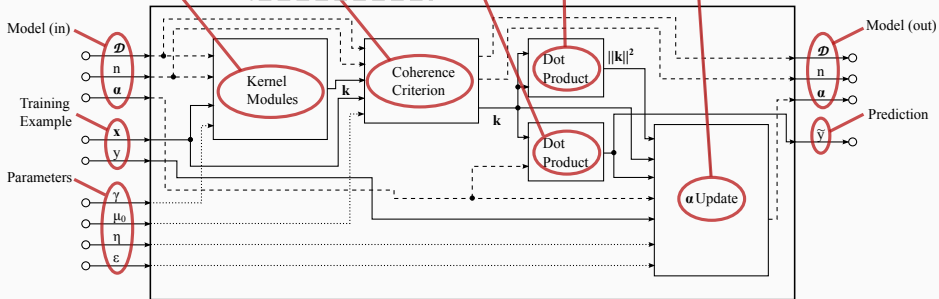
1. Create a dataflow graph (with no loops!) from a non-recursive section of the KNLMS algorithm, i.e. the update step. We call this module the *forward path*.
2. Map the dataflow graph to hardware.
3. Pipeline the hardware to achieve a desired throughput.
4. Connect an external scheduler to process parallel jobs.

A benefit of designing the hardware this way, is that the pipelining can be achieved using high level tools, such as Vivado HLS or Altera DSP Builder.

# HIGH LEVEL DESCRIPTION

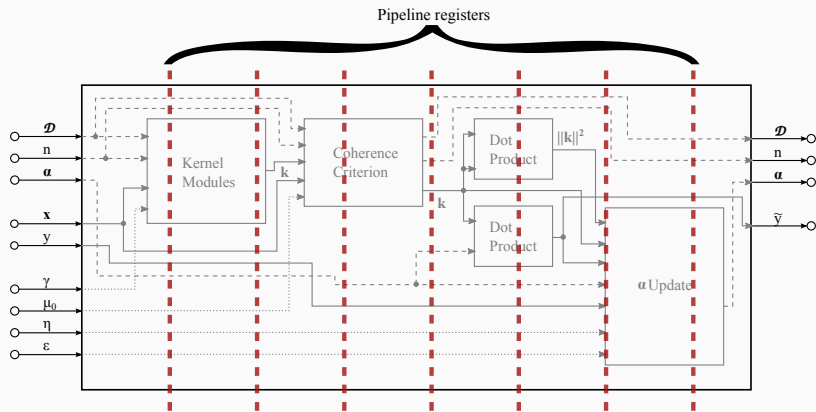
The core calculates a non-recursive portion of the KNLMS algorithm.

$$\kappa(\mathbf{x}, \tilde{\mathbf{x}}_t) = e^{-\gamma \|\mathbf{x} - \tilde{\mathbf{x}}_t\|}, \quad \mu = \max(|\mathbf{k}|)$$
$$\forall i \in \{1, \dots, N\} \quad \mathcal{D}_t = \begin{cases} [\mathcal{D}_{t-1}; \mathbf{x}] & \text{if } \mu < \mu_0 \\ \mathcal{D}_{t-1} & \text{otherwise} \end{cases}$$
$$\tilde{\mathbf{y}} = \mathbf{k}^T \boldsymbol{\alpha} \quad \|\mathbf{k}\|^2 = \mathbf{k}^T \mathbf{k} \quad \boldsymbol{\alpha}_t = \boldsymbol{\alpha}_{t-1} + \frac{\eta}{\epsilon + \|\mathbf{k}\|^2} (y_t - \tilde{y}_t) \mathbf{k}$$



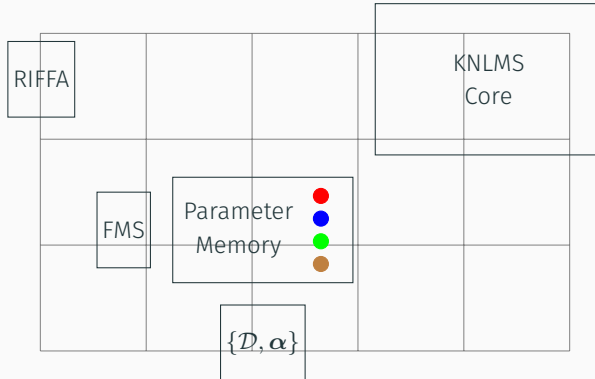
# PIPELINING

The core can be pipelined arbitrarily to achieve virtually any desired clock speed.

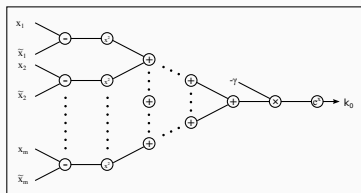




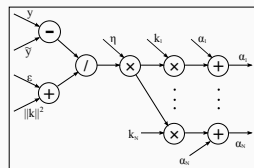
# THE PCI-SYSTEM



The complete PCI-System utilises a FSM to iterate through parameters in order to fill the KNLMS core pipeline.



**Figure 3:** Dataflow diagram of a kernel module.



**Figure 4:** Dataflow diagram of the  $\alpha$  update module.

- The design was made using Vivado HLS.
- Two variants of the core were made, one which used single precision floating point arithmetic (float), the other used a combination of fixed and single precision floating point (fused).
- A PCI-based system implementation was created using RIFFA 2.0 [Jacobsen et al., 2012]. The system implementation includes optimizations specific to a parameter search problem.
- The design was compared to: KAFBOX [Van Vaerenbergh, 2012] (MATLAB), an optimised C implementation, a naïve Vivado HLS and a previous microcoded vector processor implementation [Pang et al., 2013].

**Table 2:** Comparison of online kernel method implementations. Note that  $N = 16$  and  $M = 8$ .

Implementation	Algorithm	DSPs	Freq (MHz)	Time (ns)	Slowdown rel. to Float
Naïve	KNLMS	12	96.7	7,829	2,462
Float	KNLMS	1267	314	3.18	1
CPU (C)	KNLMS	-	3,600	940	296
CPU (KAFBOX)	KNLMS	-	3,600	73,655	23,162
Pang et al. [2013]	SW-KRLS	30	237	9,000	2,830

**Table 3:** Comparison of online kernel method implementations. Note that  $N = 16$  and  $M = 8$ .

Implementation	Algorithm	DSPs	Freq (MHz)	Time (ns)	Slowdown rel. to System
Naïve	KNLMS	12	96.7	7,829	576
Float	KNLMS	1267	314	3.18	0.23
System	KNLMS	691	250	13.6	1
CPU (C)	KNLMS	-	3,600	940	69
CPU (KAFBOX)	KNLMS	-	3,600	73,655	1703
Pang et al. [2013]	SW-KRLS	30	237	9,000	661

**Table 4:** Area utilisation of different designs obtained from synthesis. Note the computational complexity of KNLMS is  $\mathcal{O}(mn)$ .

Type	M	N	LUTs	DSPs	L	$F_{max}$
Float	2	16	77K	595	185	385
	4	16	109K	819	196	385
	16	16	307K	2163	218	385
	8	2	23K	161	162	385
	8	4	46K	319	177	385
	8	8	95K	635	192	385
	8	16	173K	1267	207	385
Fused	2	16	102K	494	161	303
	4	16	119K	595	163	303
	16	16	440K	1171	175	303
	8	2	33K	101	131	303
	8	4	64K	199	143	303
	8	8	130K	395	155	303
	8	16	247K	787	167	303