

# Braiding: a Scheme for Resolving Hazards in NORMA

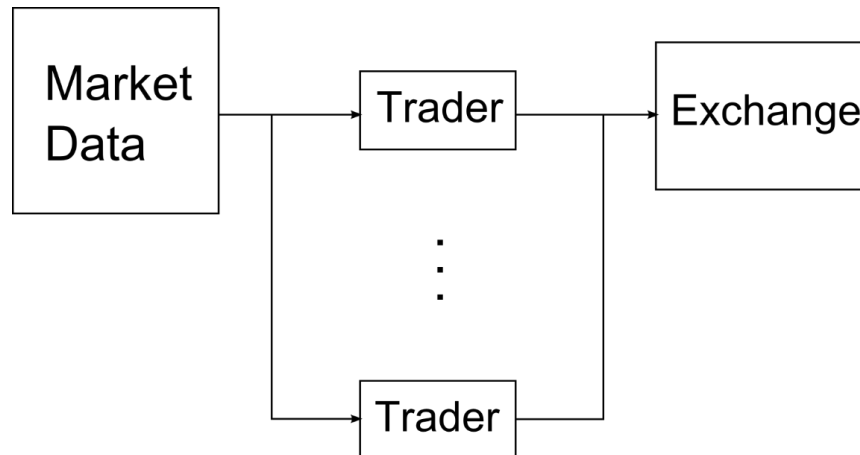
Stephen Tridgell, Duncan J.M. Moss, Nicholas J. Fraser and Philip H.W. Leong  
School of Electrical and Information Engineering,  
The University of Sydney



THE UNIVERSITY OF  
SYDNEY

How to beat other people to the money (latency)

- › Low latency trading looks to trade in transient situations where market equilibrium disturbed
  - 1ms reduction in latency can translate to \$100M per year



- › Latency also important to: prevent blackouts (cascading faults), turn off machine before it damages itself, etc

# Latency infrastructure already available

## Exablaze Low-Latency Products



ExaLINK Fusion 48 SFP+ port layer  
2 switch for replicating data typical 5  
ns fanout, 95 ns aggregation, 110 ns  
layer 2 switch

Xilinx Ultrascale FPGA, QDR SRAM,  
ARM processor



ExaNIC X10 typical raw frame  
latency 60 bytes 780 ns

**What we can't do: ML with this  
type of latency**

Source: [exablaze.com](http://exablaze.com)

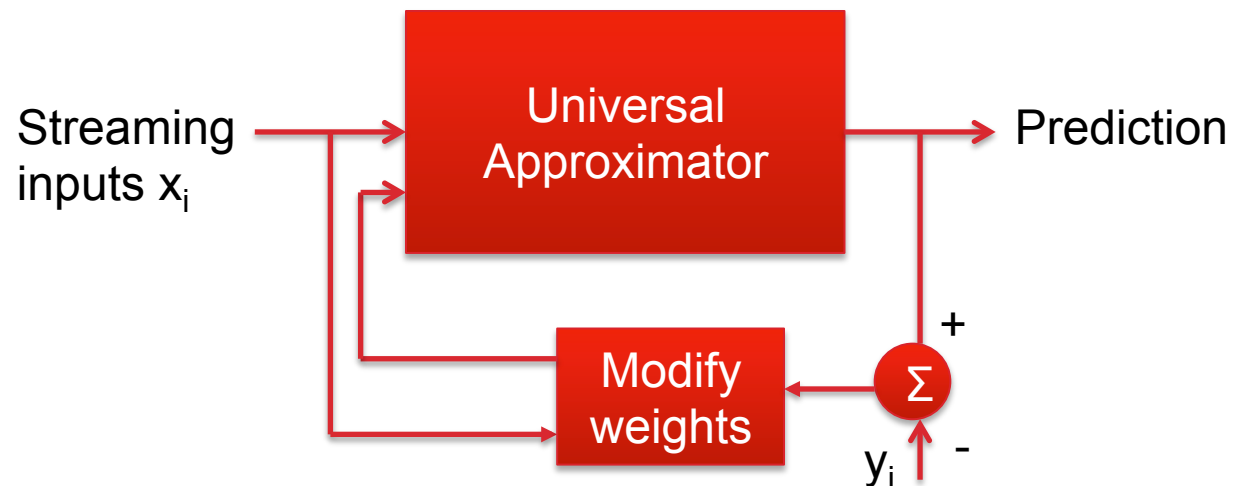
Examples are KLMS and KRLS

› Traditional ML algorithms batch based

- Several passes through data
- Requires storage of input data
- Not suitable for massive datasets

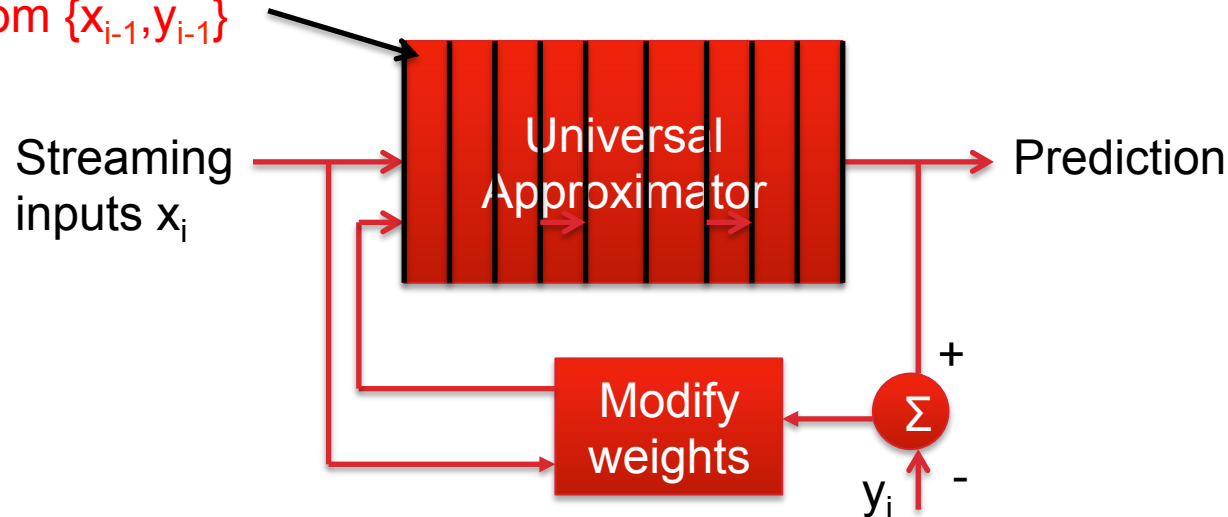
› Our approach: online algorithms

- Incremental, inexpensive state update based on new data
- Single pass through the data
- Can be high throughput, low latency



## Dependency Problem

Cannot process  
 $x_i$  until we update weights  
from  $\{x_{i-1}, y_{i-1}\}$

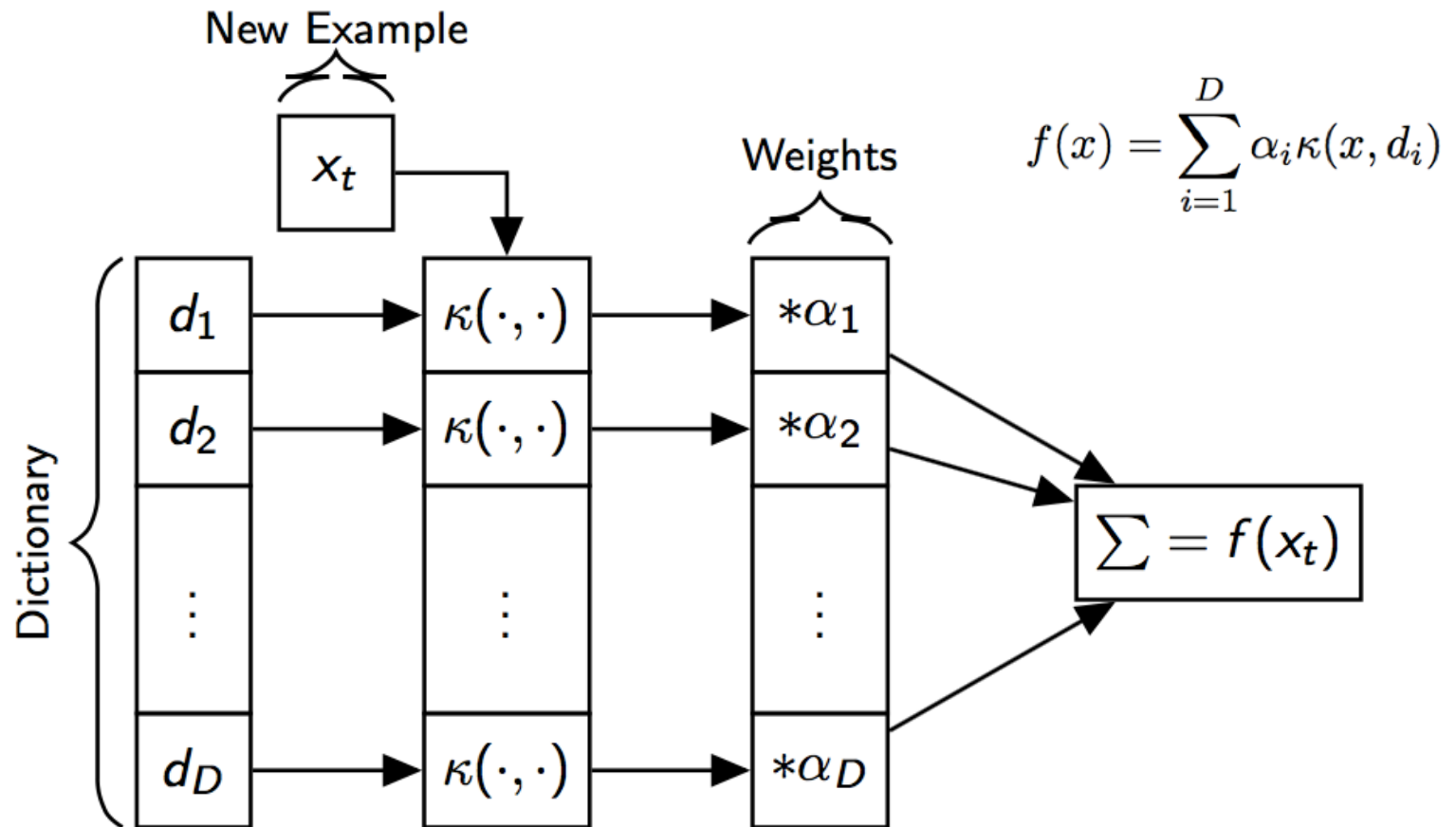


### > This work

- Implementation of an online machine learning scheme: NORMA
- Resolve read-after-write dependencies through braiding



# Datapath for NORMA



## Naive Online regularised Risk Minimization Algorithm

- › Finds Dictionary  $d_i$ , and  $\alpha_i$  (weights)

$$f(x) = \sum_{i=1}^D \alpha_i \kappa(x, d_i)$$

- › Minimise predictive error ( $R_{inst,\lambda}$ ) by taking a step in direction of gradient

$$f_{t+1} = f_t - \eta_t \partial_f R_{inst,\lambda}[f, x_{t+1}, y_{t+1}] \Big|_{f=f_t}$$

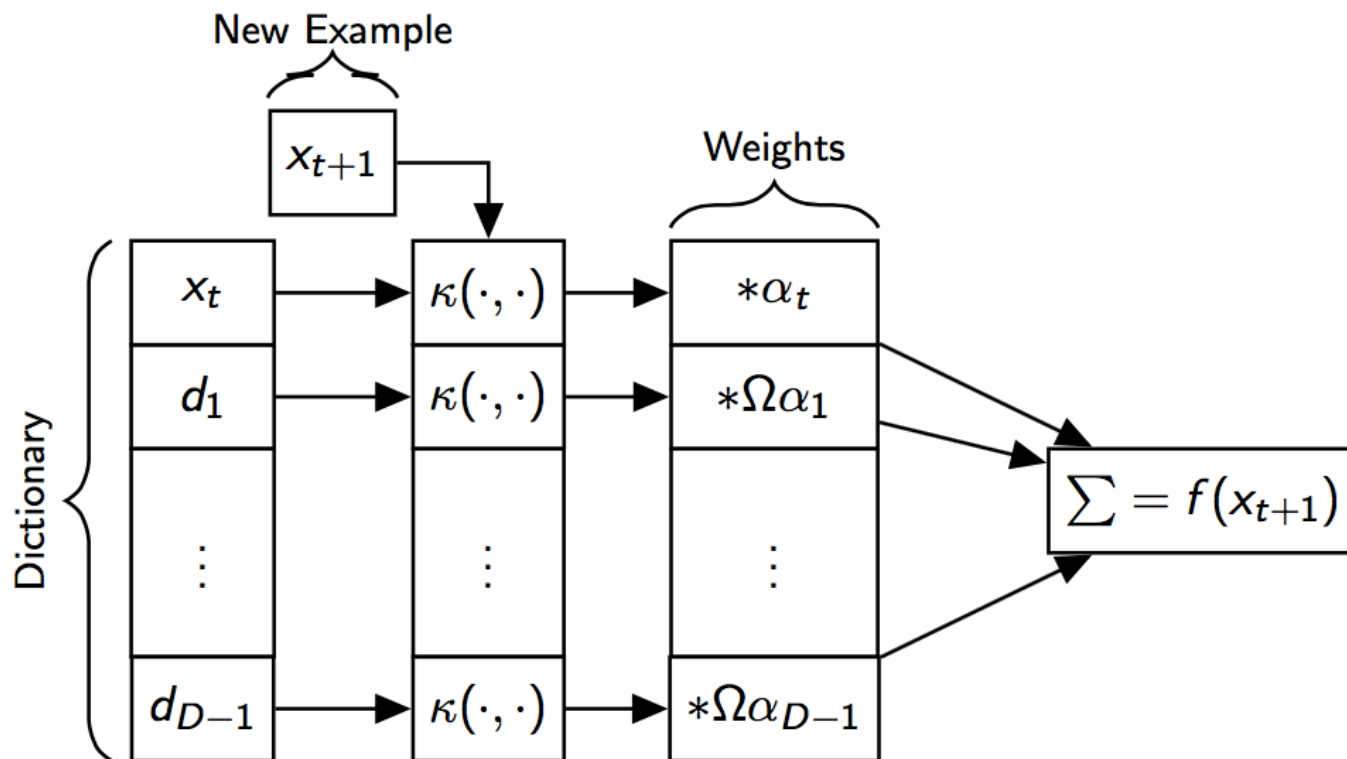
- › Can be used for classification, regression, novelty detection
- › Update for novelty detection

$$(\alpha_i, \alpha_t, \rho) = \begin{cases} (\Omega \alpha_i, 0, \rho + \eta \nu) & \text{if } f(x_t) \geq \rho \quad \text{Add } \mathbf{x}_{t+1} \text{ to dictionary} \\ (\Omega \alpha_i, \eta, \rho - \eta(1 - \nu)) & \text{otherwise} \end{cases}$$



## NORMA Update (Case 1)

$$(\alpha_i, \alpha_t, \rho) = \begin{cases} (\Omega\alpha_i, 0, \rho + \eta\nu) & \text{if } f(x_t) \geq \rho \quad (\text{Add } x_t \text{ to dictionary}) \\ (\Omega\alpha_i, \eta, \rho - \eta(1 - \nu)) & \text{otherwise} \end{cases}$$

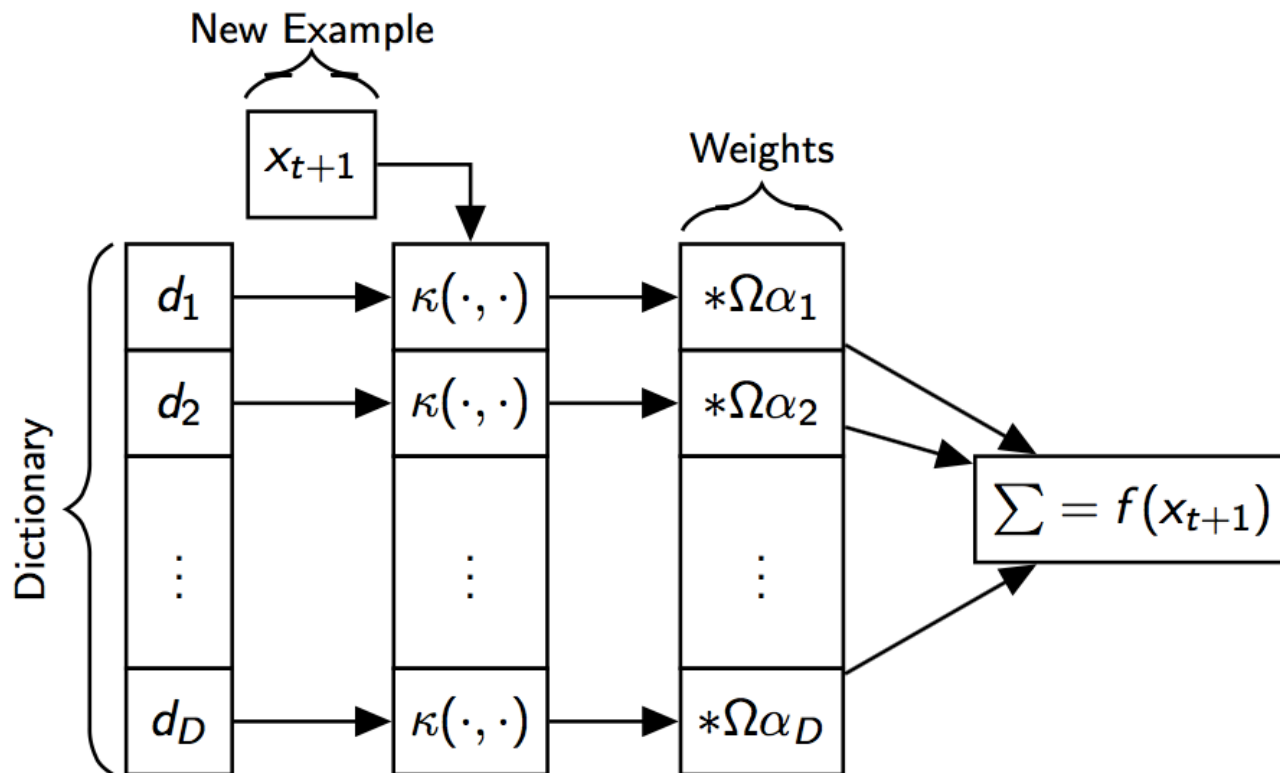






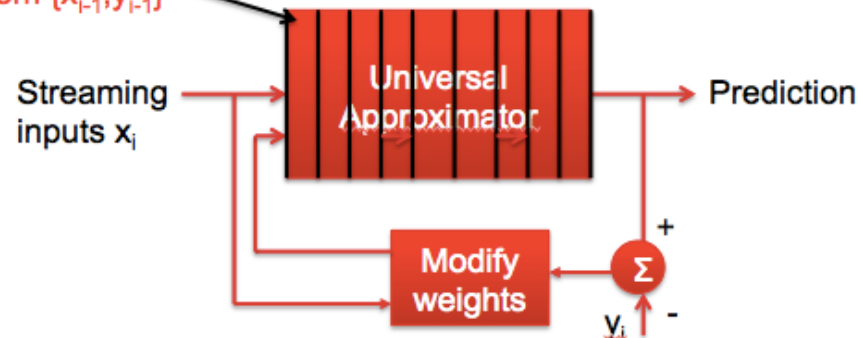
## NORMA Update (Case 2)

$$(\alpha_i, \alpha_t, \rho) = \begin{cases} (\Omega\alpha_i, 0, \rho + \eta\nu) & \text{if } f(x_t) \geq \rho \quad (\text{Add } x_t \text{ to dictionary}) \\ (\Omega\alpha_i, \eta, \rho - \eta(1 - \nu)) & \text{otherwise} \end{cases}$$

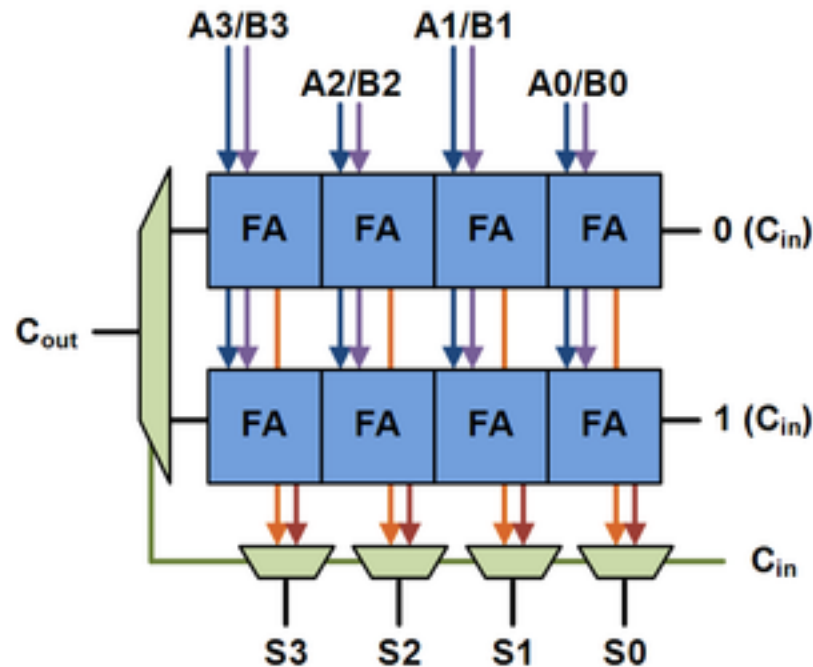


- › NORMA is a sliding window algorithm
  - If new dictionary entry added  $[d_1, \dots, d_D] \rightarrow [x_t, d_1, \dots, d_{D-1}]$
  - Weight update is just a decay  $\alpha_i \rightarrow \Omega \alpha_i$
  - Update cost is small compared to computing  $f(x_t)$
- › **Is this really true?**

Cannot process  
 $x_i$  until we update weights  
from  $\{x_{i-1}, y_{i-1}\}$



- › Recall carry select adder
  - implement both cases in parallel and select output



Source: Wikipedia



$$f(x_{t+1}) = \sum_{i=1}^D \alpha_i \kappa(x_{t+1}, d_i)$$

Use the previous dictionary for  $x_t$  denoted  $\hat{d}_i$

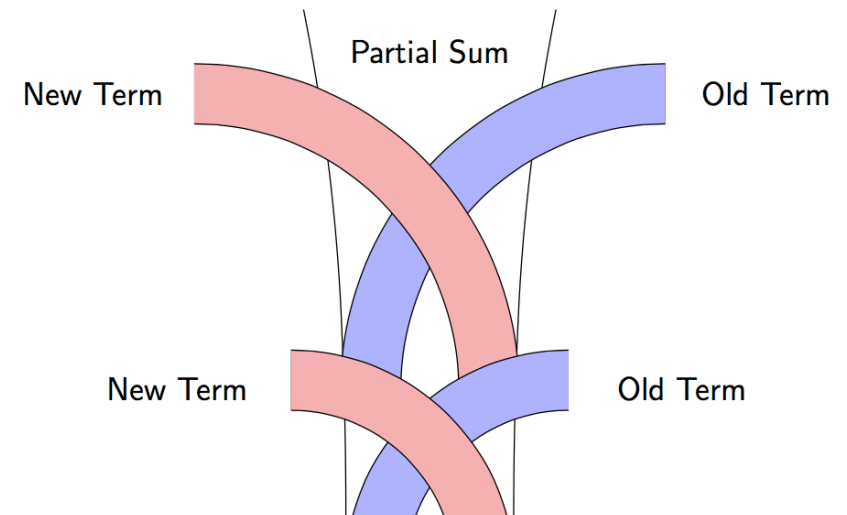
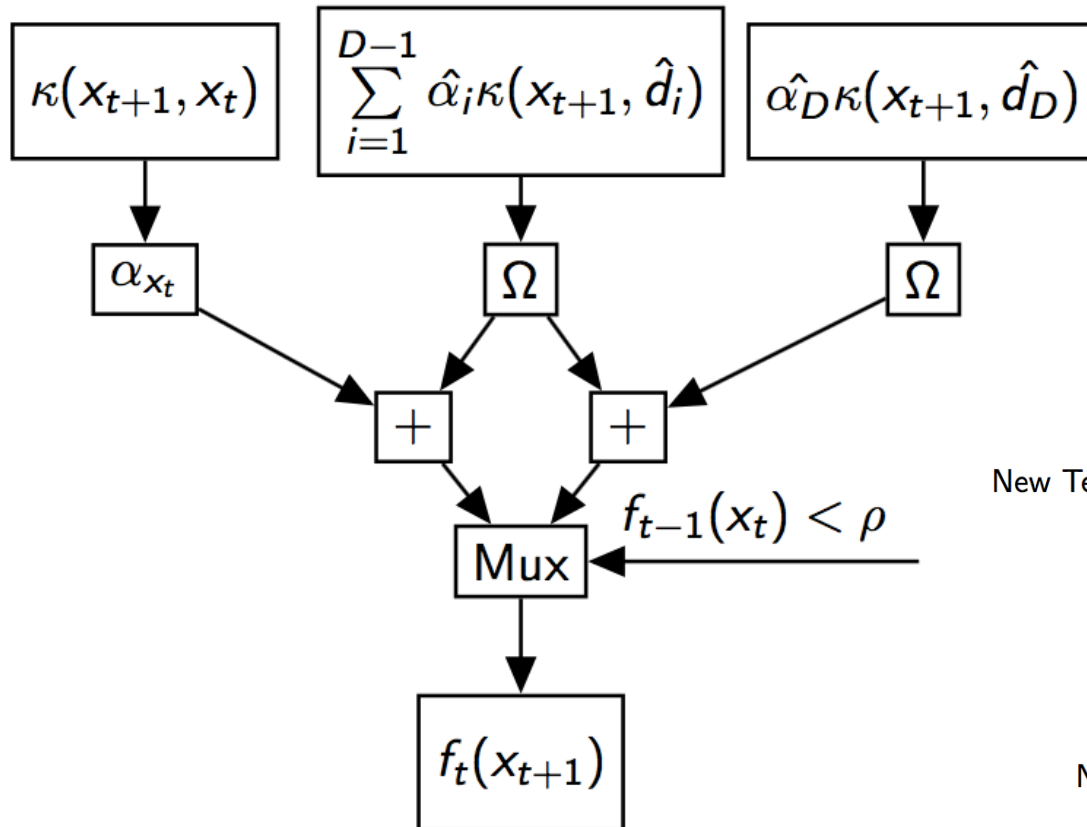
$$f(x_{t+1}) = \sum_{i=1}^{D-1} \Omega \hat{\alpha}_i \kappa(x_{t+1}, \hat{d}_i) + \text{something}$$

if  $x_t$  is added then this term =  $\alpha_{x_t} \kappa(x_{t+1}, x_t)$

if  $x_t$  is not added then this term =  $\Omega \hat{\alpha}_D \kappa(x_{t+1}, \hat{d}_D)$



# Braiding Datapath



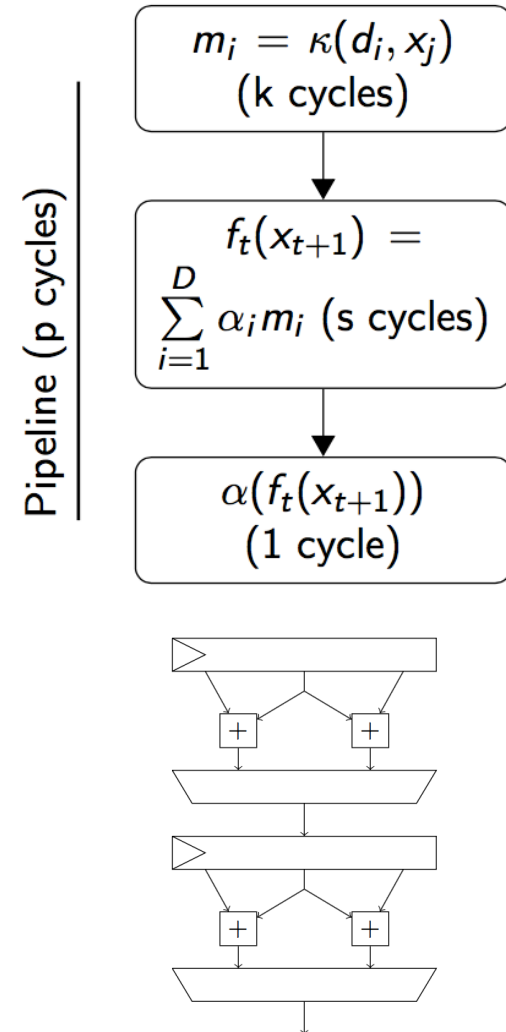


# Generalised to p cycles

$$f_t(x_{t+1}) = \sum_{i=1}^{D-p} \Omega^p \hat{\alpha}_i \kappa(x_{t+1}, \hat{d}_i)$$

$$q \left\{ \begin{array}{l} + \left\{ \begin{array}{l} 0 \text{ if } x_{t+1-p} \text{ is not added} \\ \Omega^{p-1} \alpha_{x_{t+1-p}} \kappa(x_{t+1}, x_{t+1-p}) \text{ otherwise} \end{array} \right. \\ + \left\{ \begin{array}{l} 0 \text{ if } x_{t+2-p} \text{ is not added} \\ \Omega^{p-2} \alpha_{x_{t+2-p}} \kappa(x_{t+1}, x_{t+2-p}) \text{ otherwise} \end{array} \right. \\ \vdots \\ + \left\{ \begin{array}{l} 0 \text{ if } x_t \text{ is not added} \\ \alpha_{x_t} \kappa(x_{t+1}, x_t) \text{ otherwise} \end{array} \right. \end{array} \right.$$

$$+ \sum_{i=D-p+1}^{D-q} \Omega^p \hat{\alpha}_i \kappa(x_{t+1}, \hat{d}_i)$$



- › Implemented in Chisel
- › On XC7VX485T- 2FFG1761C achieves ~133 MHz
- › Area  $O(FDB^2)$  ( $F$ =dimensionality of input vector), time complexity  $O(FD)$
- › Speedup 500x compared with single core CPU i7-4510U (8.10 fixed)

<b>F=8, D=</b>	<b>16</b>	<b>32</b>	<b>64</b>	<b>128</b>	<b>200</b>
<b>Frequency (MHz)</b>	133	138	137	131	127
<b>DSPs (/2,800)</b>	309	514	911	1,679	2,556
<b>Slices (/759,000)</b>	4,615	8,194	14,663	29,113	46,443
<b>Latency (cycles)</b>	10	11	12	12	13
<b>Speedup (×)</b>	47	91	178	344	509
<b>Latency reduction (×)</b>	4.69	8.30	14.9	28.7	39.2

# Comparison of Architectures

- › Core with input vector  $F=8$  and dictionary size  $D=16$

Design	Precision	Freq MHz	Latency Cycles	T.put Cycles	Latency nS	T.put nS
<b>Vector KNLMS</b>	Single	282	479	479	1,699	1,699
<b>Pipelined KNLMS</b>	Single	314	207	1	659	3.2
<b>Braided NORMA</b>	8.10	113	10	1	89	8.8



- › Braiding: rearrangement of a sliding window algorithm for hardware implementations
  - NORMA used but other ML algorithms possible
- › Compared with pipelined KNLMS,
  - 20x lower latency at 1/3 of the throughput
- › Open source (GPLv2): [github.com/da-steve101/chisel-pipelined-olk](https://github.com/da-steve101/chisel-pipelined-olk)