

ESE532: System-on-a-Chip Architecture

Day 25: November 25, 2019
Reduce



Today

- Reduce
- Associative Operations
- Model
- Latency Bound Implications and Implementations
- Parallel Prefix
- Broad Application

Message

- Aggregation is a common need that is not strictly data parallel
- ...but admits to parallel computation

Reduce

- Reduce – combining a collection of data into a single value
 - Converting a vector into a scalar

Sum Reduce

- Simplest and most common
 - Add up all the values in a vector or array

```
int sum=0;
for (int i=0;i<N; i++)
    sum+=a[i];
```

Sum Reduce

- What's II?

```
int sum=0;
for (int i=0;i<N; i++)
    sum+=a[i];
```

Sum Reduce

- What's latency bound?
 - Assuming associativity holds for addition

```
int sum=0;
for (int i=0;i<N; i++)
    sum+=a[i];
```

Associative Operations

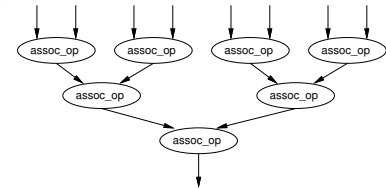
- Associativity means can group together operations in any way
- Normal sequential:
 $((a[0]+a[1])+a[2])+a[4])+...$
- Associative regroup:
 $(a[0]+(((a[1]+(a[2]+a[3]))+a[4])+(...)))$

Associative Operations

- Associativity means can group together operations in any way
- Normal sequential:
 $((a[0]+a[1])+a[2])+a[4])+...$
- Regroup parallelism:
 $((a[0]+a[1])+(a[2]+a[3]))+((a[4]+a[5])+(a[6]+a[7]))$

Associative Tree Reduce

- Add pairs – cut numbers in half
- Repeat adding pairs until single value
- How deep?



Latency Bounds

- Associative reduces typically contribute log terms to latency bounds
 - ...as you've seen on many previous midterms and finals

Sum Reduce

- Data Parallel?

```
int sum=0;
for (int i=0;i<N; i++)
    sum+=a[i];
```

Sum Reduce

- How exploit 4 cores to compute?
 - (assume a very large, like 1 million)

```
int sum=0;
for (int i=0;i<N; i++)
    sum+=a[i];
```

Model: Data Parallel+Reduce

- Data Parallel + Reduce
 - Very common to perform a data parallel operation then a reduce on results

- Example: dot product

```
int sum=0;
for (int i=0;i<N; i++)
    sum+=a[i]*b[i];
```

Dot Product

- Latency bound for dot product
 - Assume 1 cycle add, 3 cycle multiply

- Example: dot product

```
int sum=0;
for (int i=0;i<N; i++)
    sum+=a[i]*b[i];
```

Model: Data Parallel+Reduce

- Data Parallel + Reduce
 - Very common to perform a data parallel operation then a reduce on results

- General form

```
int res=0;
for (int i=0;i<N; i++)
    res=assoc_op(res,f(a[i],b[i], ...))
```

What else Associative?

- Beyond addition, what other associative operations do we often see as reductions?

Associative Operations

- Add
- Multiply
- Max
- Min
- AND
- OR
- Max/min
 - And keep associated values
- Find First

Optimization Loop

```
int minval=f(0);
int min=0;
for (i=1;i<N;i++) {
    int val=f(i);
    if (val<minval) {
        minval=val; min=i;
    }
}
```

Penn ESE532 Fall 2019 -- DeHon

19

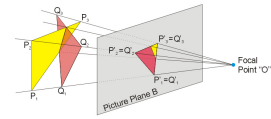
Rendering Decomposed

Day 15

• Pipeline of

– Projection

- Where do the points of this triangle end up in the viewed image?
- Matrix-multiplication to translate points



– Rasterization

- Turn into pixels
- Fill pixels for triangle



Figures from:
https://commons.wikimedia.org/wiki/File:Perspective_Projection_Principle.jpg
https://en.wikipedia.org/wiki/Rasterization#/media/File:Raster_graphic_fish_20x23squares_sd_tv-example.png

– Z-buffer

- Keep only the ones on top (not hidden)
 - 2D image + Z-depth – keep smallest

Penn ESE532 Fall 2019 -- DeHon

20

Z-Buffering

- Storing into Z-buffer is an associative reduce operation
 - Min reduce (keep nearest pixel) on depth with an associated value
- Parallel strategy
 - Split triangles into sets
 - Project, rasterize, Z-buffer in parallel
 - Assoc. reduce Z-buffer pixels across parallel Z-buffers

Penn ESE532 Fall 2019 -- DeHon

21

Data Parallel+Preduce

IMPLEMENTATIONS

Penn ESE532 Fall 2019 -- DeHon

22

Threaded: Data Parallel+Reduce

- Break into P threads
 - 0 to N/P-1, N/P to 2N/P-1, ...
- Run fraction of data and reduce on each
- Then bring results together to sum
 - P small, on one processor
 - P large, as tree

Penn ESE532 Fall 2019 -- DeHon

23

Model: Data Parallel+Reduce

• What's cycle → what's II?

• General form

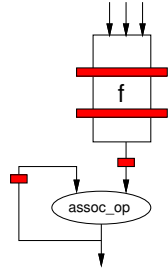
```
int res=0;
for (int i=0;i<N; i++)
    res=assoc_op(res,f(a[i],b[i], ...))
```

Penn ESE532 Fall 2019 -- DeHon

24

Pipeline: Data Parallel + Reduce

- Pipeline f
- Cycle on assoc_op

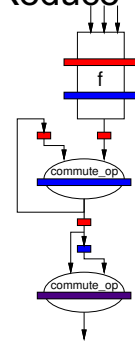


Penn ESE532 Fall 2019 -- DeHon

25

Pipeline: Data Parallel + Reduce

- Pipeline f
- Cycle on assoc_op
- Avoid cycle, $II=1$ for commutative
 - Run interleaved
 - C-slow ($C=II$)
 - Combine at end
- Commutative operation
 - Can reorder operands

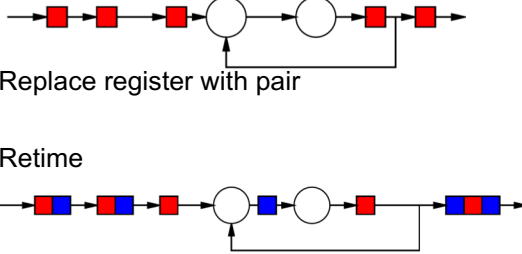


Penn ESE532 Fall 2019 -- DeHon

26

2-Slow Simple Cycle Day 7

- Replace register with pair
- Retime
- Observe independence of red/blue computations

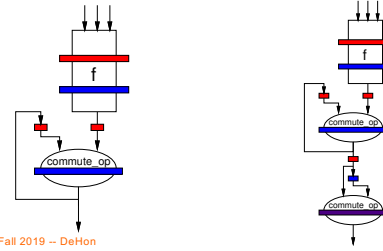


Penn ESE532 Fall 2019 -- DeHon

27

$II=2$ Commutative Pipelined Reduce

- 2-slow transformation
- Combine independent (red, blue) reduces

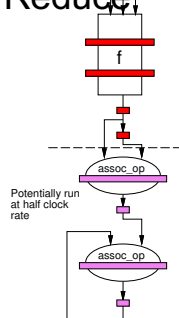


Penn ESE532 Fall 2019 -- DeHon

28

Pipeline: Data Parallel + Reduce

- Pipeline f
- Cycle on assoc_op
- Avoid cycle, $II=1$ for associative
 - Gather up II values
 - Run through pipelined assoc. reduce tree
 - Drop into assoc_op cycle every II cycles



Penn ESE532 Fall 2019 -- DeHon

29

Model: Data Parallel+Reduce

- **Conclude:** associative reduce can achieve II of 1
- General form


```
int res=0;
for (int i=0;i<N; i++)
    res=assoc_op(res,f(a[i],b[i], ...))
```

Penn ESE532 Fall 2019 -- DeHon

30

Vector:

Data Parallel + Reduce

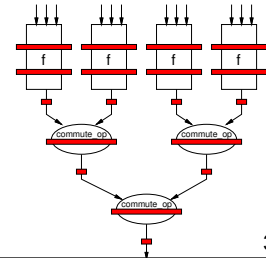
- Some vector/SIMD machines will have dedicated reduce hardware
 - E.g. vector-add operator
 - NEON
 - Not have vector reduce
 - Does have VPADAL
 - Use VL adds for course-grained reduce
- ```
for (i=0; i<N; i+=VL) {
 avl=a[i]...a[i+VL-1]
 VADD(res,avl, res);
}
```
- Use VPADAL to complete

Penn ESE532 Fall 2019 -- DeHon

31

### Unrolled Pipeline: Data Parallel + Reduce

- Unroll computation
- Perform f ops in parallel pipelines
- Pipelined tree reduce



Penn ESE532 Fall 2019 -- DeHon

32

## PARALLEL PREFIX

Penn ESE532 Fall 2019 -- DeHon

33

## What if want Prefix?

### Sum Reduce

```
int sum=0;
for (int i=0; i<N; i++)
 sum+=a[i];
```

### Sum Prefix

```
int sum[N];
sum[0]=a[0];
for (int i=1; i<N; i++)
 sum[i]=a[i]+sum[i-1];
```

Penn ESE532 Fall 2019 -- DeHon

34

## Prefix

- Aggregate (vector) output where item  $i$  is the reduce of the input vector 0 through  $i$
- ```
prefix[0]=a[0];  
for (int i=1; i<N; i++)  
    prefix[i]=op(prefix[i-1], f(a[i]...));
```

Penn ESE532 Fall 2019 -- DeHon

35

Latency Bound

- What's the latency bound for the prefix when op is associative?
 - Assume op is 1 cycle
 - How $cycles(op) > 1$ change?
- ```
prefix[0]=a[0];
for (int i=1; i<N; i++)
 prefix[i]=op(prefix[i-1], f(a[i]...));
```

Penn ESE532 Fall 2019 -- DeHon

36

## Resources?

- How much hardware to achieve latency bound?

```
prefix[0]=a[0];
for (int i=1;i<N; i++)
 prefix[i]=op(prefix[i-1],f(a[i]...));
```

## Latency Bound?

- What's the latency bound for this operation?

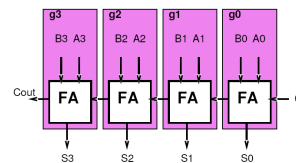
```
boolean a[i],b[i],s[i]
for (i=0;i<N;i++) {
 cn=(a[i]&&b[i])||
 (a[i]&&c)||
 (b[i]&&c);
 s[i]=a[i] ^ b[i] ^ c;
 c=cn;
}
```

## Associative

- Is the carry operation in addition an associative operation?
- Operation:
  - MAJ=majority =  $A \& B \ || \ B \& C \ || \ A \& B$

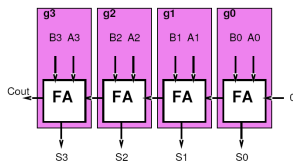
## Carry Computation

- Think about each adder bit as a computing a function on the carry in
  - $C[i]=g(c[i-1])$
  - Particular function  $f$  will depend on  $a[i], b[i]$
  - $g=f(a,b)$



## Functions

- What are the functions  $g(c[i-1])$ ?
  - $g(c)=\text{carry}(a=0,b=0,c)$
  - $g(c)=\text{carry}(a=1,b=0,c)$
  - $g(c)=\text{carry}(a=0,b=1,c)$
  - $g(c)=\text{carry}(a=1,b=1,c)$

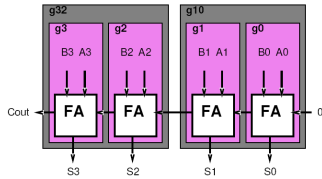


## Functions

- What are the functions  $g(c[i-1])$ ?
  - $g(x)=1$  **Generate**
    - $a[i]=b[i]=1$
  - $g(x)=x$  **Propagate**
    - $a[i] \text{ xor } b[i]=1$
  - $g(x)=0$  **Squash**
    - $a[i]=b[i]=0$

## Combining

- Want to combine functions
  - Compute  $c[i]=g_i(g_{i-1}(c[i-2]))$
  - Compute compose of two functions
- What functions will the compose of two of these functions be?
  - Same as before
    - Propagate, generate, squash



## Compose Rules (LSB MSB)

- GG
- GP
- GS
- PG
- PP
- PS
- SG
- SP
- SS

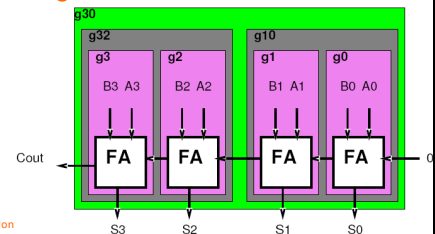
[work on board]

## Compose Rules (LSB MSB)

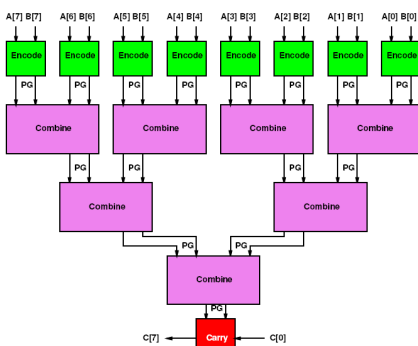
- GG = G
- GP = G
- GS = S
- PG = G
- PP = P
- PS = S
- SG = G
- SP = S
- SS = S

## Combining

- Do it again...
- Combine  $g[i-3, i-2]$  and  $g[i-1, i]$
- What do we get?

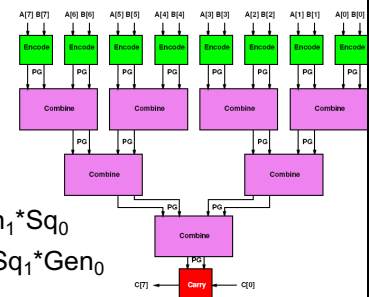


## Associative Reduce Tree



## Reduce Tree

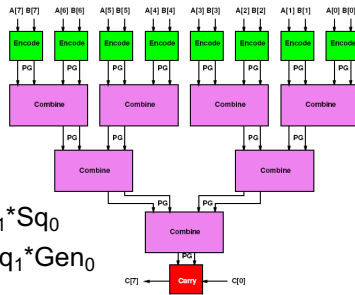
- $Sq = A * B$
- $Gen = A * B$
- $Sq_{out} = Sq_1 + Gen_1 * Sq_0$
- $Gen_{out} = Gen_1 + Sq_1 * Gen_0$





## Reduce Tree

- $Sq = /A*/B$
- $Gen = A*B$
- $Sq_{out} = Sq_1 + /Gen_1 * Sq_0$
- $Gen_{out} = Gen_1 + /Sq_1 * Gen_0$
- Delay and Area? (work next few slides)



Penn ESE532 Fall 2019 -- DeHon

49

## Reduce Tree

- $Sq = /A*/B$
- $Gen = A*B$
- $Sq_{out} = Sq_1 + /Gen_1 * Sq_0$
- $Gen_{out} = Gen_1 + /Sq_1 * Gen_0$
- $A(Encode) = 2$
- $D(Encode) = 1$
- $A(Combine) = 4$
- $D(Combine) = 2$
- $A(Carry) = 2$
- $D(Carry) = 1$

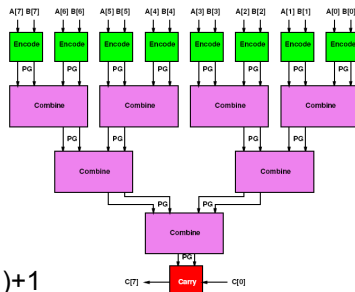
Penn ESE532 Fall 2019 -- DeHon

50

## Reduce Tree: Delay?

- $D(Encode) = 1$
- $D(Combine) = 2$
- $D(Carry) = 1$

$$\text{Delay} = 1 + 2\log_2(N) + 1$$



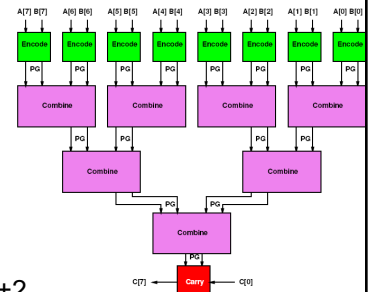
Penn ESE532 Fall 2019 -- DeHon

51

## Reduce Tree: Area?

- $A(Encode) = 2$
- $A(Combine) = 4$
- $A(Carry) = 2$

$$\text{Area} = 2N + 4(N-1) + 2$$

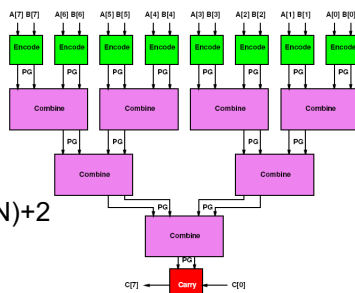


Penn ESE532 Fall 2019 -- DeHon

52

## Reduce Tree: Area & Delay

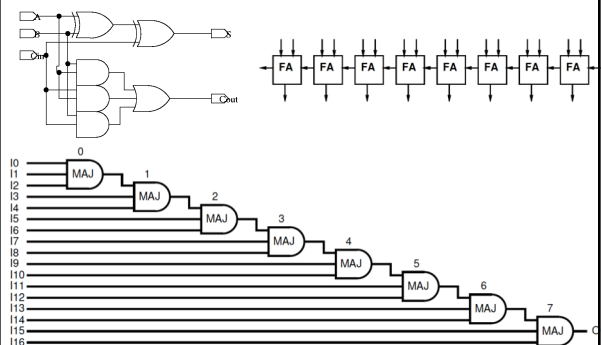
- $\text{Area}(N) = 6N - 2$
- $\text{Delay}(N) = 2\log_2(N) + 2$



Penn ESE532 Fall 2019 -- DeHon

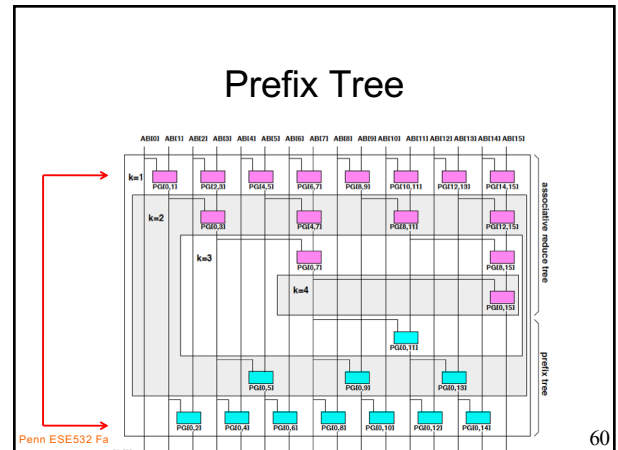
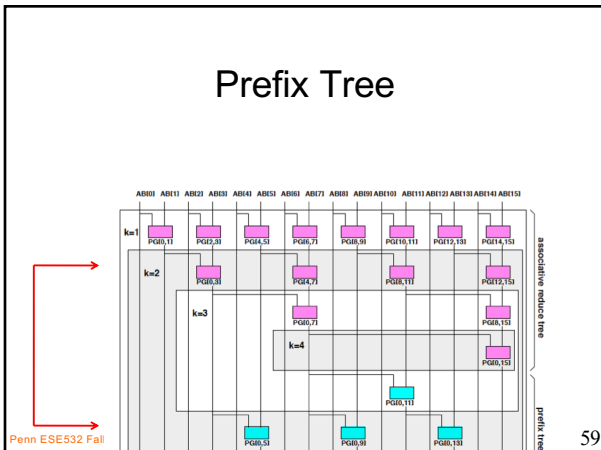
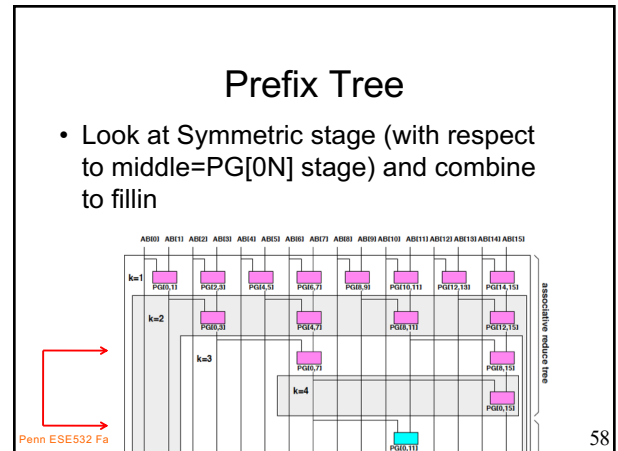
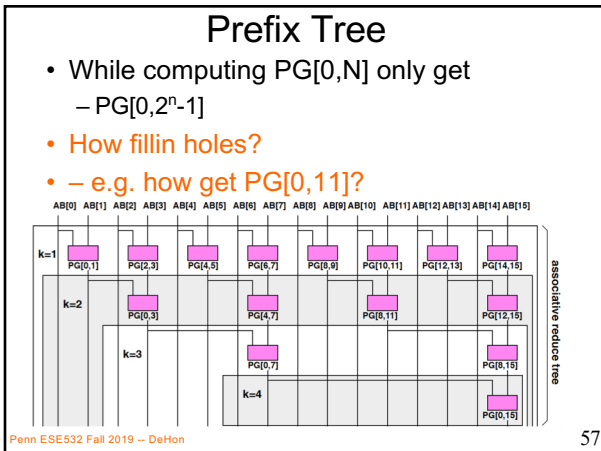
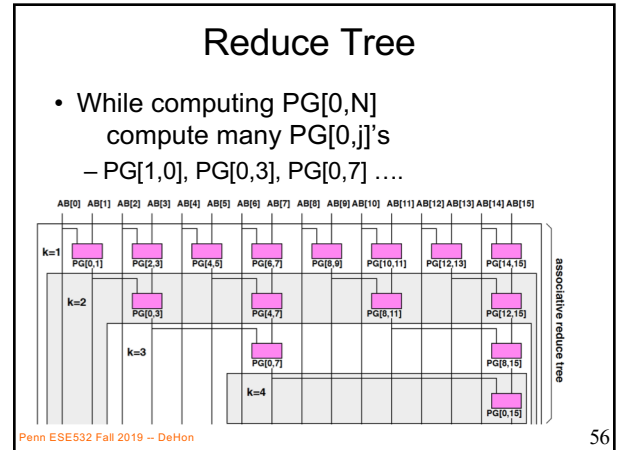
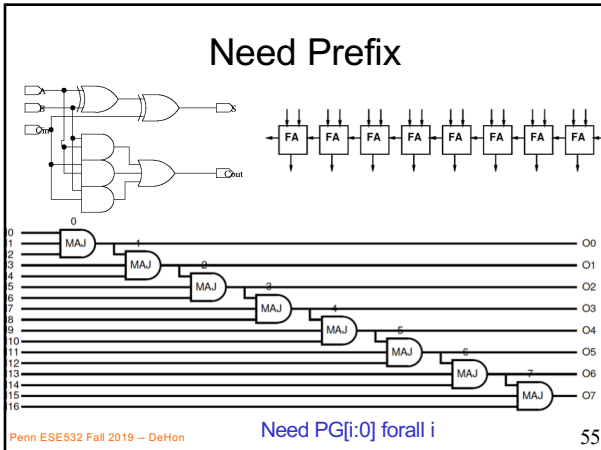
53

## Compute Carry[N]



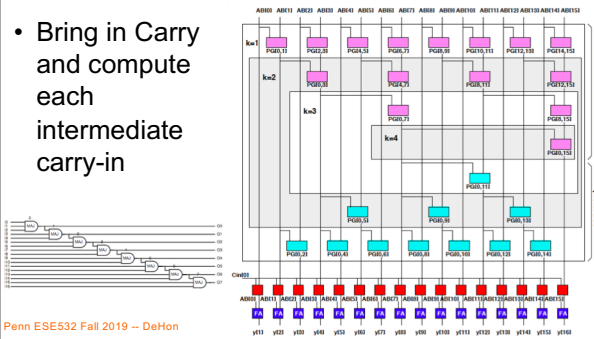
Penn ESE532 Fall 2019 -- DeHon

54



## Prefix Tree

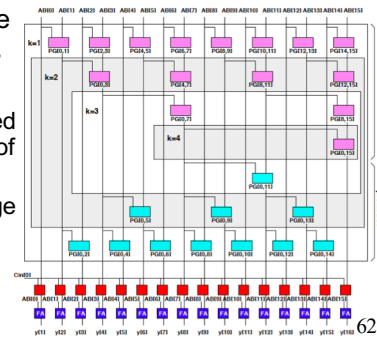
- Bring in Carry and compute each intermediate carry-in



Penn ESE532 Fall 2019 -- DeHon

## Prefix Tree

- Note: prefix-tree is same size as reduce tree
  - Always matched same number of elements in symmetric stage

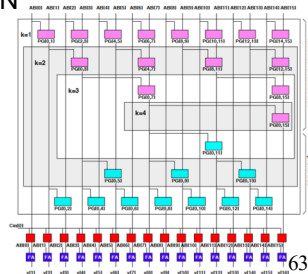


Penn ESE532 Fall 2019 -- DeHon

62

## Parallel Prefix Area and Delay?

- Roughly twice the area/delay
- Area =  $2N + 4N + 4N + 2N = 12N$
- Delay =  $4\log_2(N) + 2$
- Conclude: can add in log time with linear area.



Penn ESE532 Fall 2019 -- DeHon

63

## Parallel Prefix

- Important **Pattern**
- Applicable any time operation is *associative*
  - Or can be made assoc. as in MAJ case
- Function Composition is always associative
- Logarithmic delay
- Linear area

Penn ESE532 Fall 2019 -- DeHon

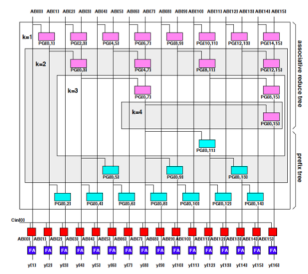
64

## Parallel Prefix Sum

```

prefix[0]=a[0];
for (int i=1;i<N; i++)
 prefix[i]=op(prefix[i-1],f(a[i]...));

```



Penn ESE532 Fall 2019 -- DeHon

65

## BROADER APPLICATION

Penn ESE532 Fall 2019 -- DeHon

66

## Cast Associative

- If you can cast it into an associative operation, you can apply
  - Associative Reduce
  - Parallel Prefix

## Examples

- Saturated Addition
  - Not associative
- Floating-Point Addition
- Finite Automata Evaluation
  
- (papers in supplemental reading)

## Big Ideas:

- Reduce from aggregate to scalar
  - is a common operation
  - not strictly data parallel
  - Associative reduce admits to parallelism
    - $\log(N)$  latency bound
    - $II=1$
    - Linear area
- Prefix when want reduce of all prefixes
  - Also  $\log(N)$  latency bound
  - Linear area

## Admin

- Wednesday is a virtual Friday
  - Will not meet
  - Happy Thanksgiving
  - Nothing due Friday (virtual or real)
- Project due **following** Friday (12/6)