

ESE532: System-on-a-Chip Architecture

Day 21: November 16, 2020
Reduce



Penn ESE532 Fall 2020 -- DeHon

Today

- Part 1
 - Reduce
 - Associative Operations
 - Model
- Part 2
 - Latency Bound Implications and Implementations
- Part 3
 - Parallel Prefix
 - Broad Application
- Bonus: Binary Arithmetic

Penn ESE532 Fall 2020 -- DeHon

2

Message

- Aggregation is a common need that is not strictly data parallel
- ...but admits to parallel computation with a slightly different pattern that is worth knowing

Penn ESE532 Fall 2020 -- DeHon

3

Reduce

- Reduce – combining a collection of data into a single value
 - Converting a vector into a scalar
 - E.g. sum elements

Penn ESE532 Fall 2020 -- DeHon

4

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];
```

Penn ESE532 Fall 2020 -- DeHon

5

Sum Reduce

- What's II?

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

Penn ESE532 Fall 2020 -- DeHon

6

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:

$$(((a[0]+a[1])+a[2])+a[3])+...$$
- Associative regroup:

$$(a[0]+(((a[1]+(a[2]+a[3]))+a[4])+(...)))$$

Associative Operations

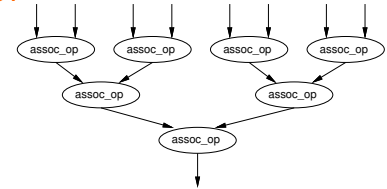
- Associativity means can group together operations in any way
- Normal sequential:

$$(((a[0]+a[1])+a[2])+a[3])+...$$
- Regroup parallelism:

$$(((a[0]+a[1])+(a[2]+a[3]))+((a[4]+a[5])+(a[6]+a[7])))$$

Associative Tree Reduce

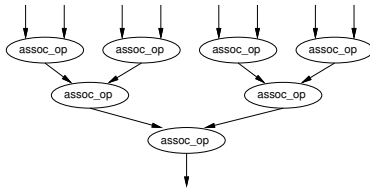
- Add pairs – cut numbers in half
- Repeat adding pairs until single value
- How deep?



Associative Tree Reduce

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

- $N \cdot (1/2)^k = 1$
- $N = 2^k$
- $k = \log_2(N)$



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
(core in DNN, Matrix-Multiply)

```
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

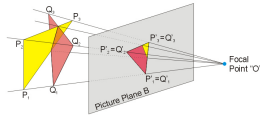
- Add
- Multiply
- Max
- Min
- AND
- OR
- Max/min
 - And keep associated values
- Find First

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;
    }
}
```

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



Figures from:
https://commons.wikimedia.org/wiki/File:Perspective_Projection_Principle.png
https://en.wikipedia.org/wiki/Rasterisation#/media/File:Raster_graphic_fish_20x23squares_sdtv-example.png

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

Part 2: Data Parallel+Reduce

IMPLEMENTATIONS

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

Model: Data Parallel+Reduce

- What's cycle \rightarrow what's II?
- General form

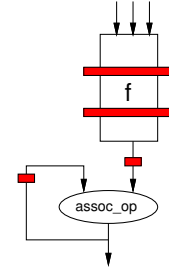

```
int res=0;
for (int i=0;i<N; i++)
  res=assoc_op(res,f(a[i],b[i], ...))
```

Penn ESE532 Fall 2020 -- DeHon

25

Pipeline: Data Parallel + Reduce

- Pipeline f
- Cycle on assoc_op

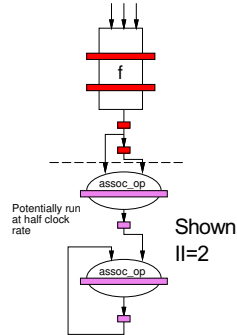


Penn ESE532 Fall 2020 -- DeHon

26

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

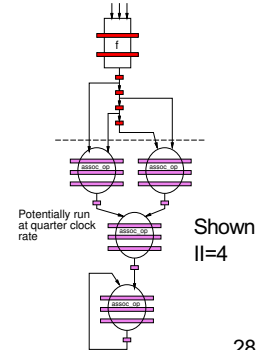


Penn ESE532 Fall 2020 -- DeHon

27

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
- Solves, but **underutilizes at 1/4 clock**



Penn ESE532 Fall 2020 -- DeHon

28

Model: Data Parallel+Reduce

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

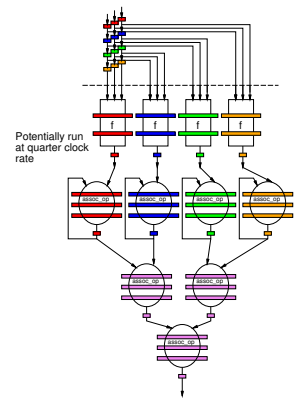

```
int res=0;
for (int i=0;i<N; i++)
  res=assoc_op(res,f(a[i],b[i], ...))
```

Penn ESE532 Fall 2020 -- DeHon

29

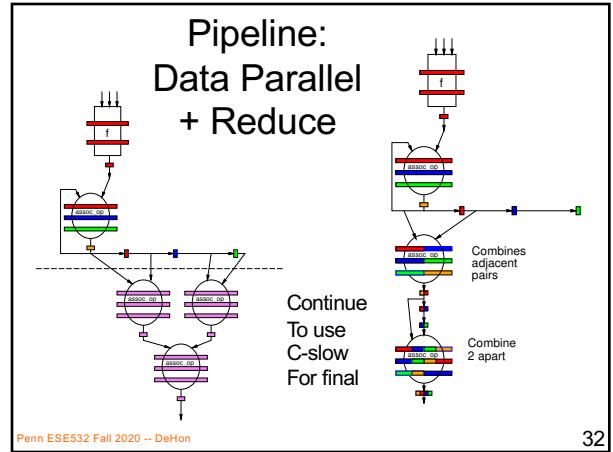
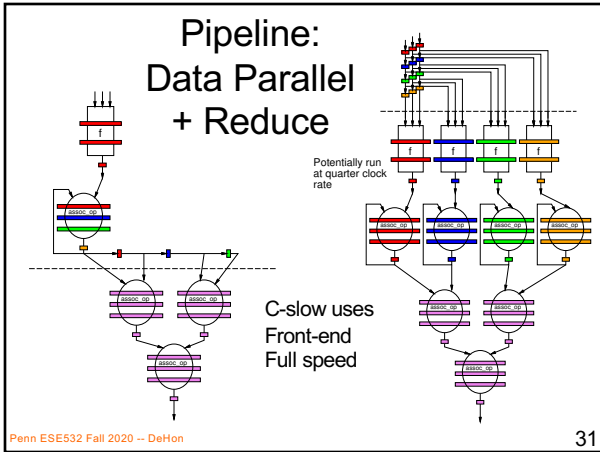
Pipeline: Data Parallel + Reduce

- Alternate, less efficient
 - ...but setting up next



Penn ESE532 Fall 2020 -- DeHon

30

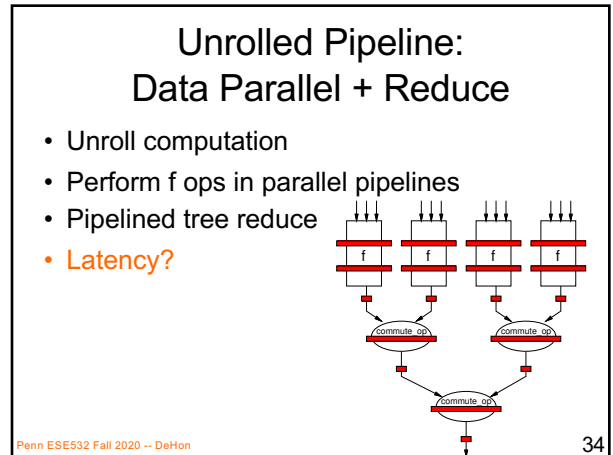


Vector: Data Parallel + Reduce

- Some vector/SIMD machines will have dedicated reduce hardware
- E.g. vector-add operator
- NEON
 - Not have vector reduce
 - Does have VPADAL
 - Pairwise adds
- 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 2020 -- DeHon 33



Implement Reduce

- Can exploit with all of our parallel implementation forms
 - Multi-thread (multi-processor)
 - SIMD/Vector
 - Pipeline
 - Spatial (unrolled)

Penn ESE532 Fall 2020 -- DeHon 35

Part 3

PARALLEL PREFIX

Penn ESE532 Fall 2020 -- DeHon 36

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 2020 -- DeHon

37

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 2020 -- DeHon

38

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 2020 -- DeHon

39

Resources?

- How much hardware to achieve latency bound?

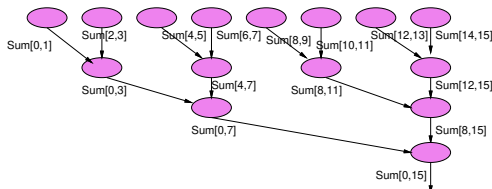
```
prefix[0]=a[0];
for (int i=1;i<N; i++)
    prefix[i]=op(prefix[i-1],f(a[i]...));
```

Penn ESE532 Fall 2020 -- DeHon

40

Reduce Tree

- While computing $Sum[0,N-1]$ compute many $Sum[0,j]$'s
 - $Sum[0,1]$, $Sum[0,3]$, $Sum[0,7]$

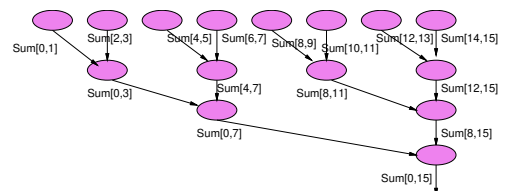


Penn ESE532 Fall 2020 -- DeHon

41

Prefix Tree

- While computing $Sum[0,N-1]$ only get
 - $PG[0,2^n-1]$
- How fillin holes?
 - e.g. how get $Sum[0,11]$?

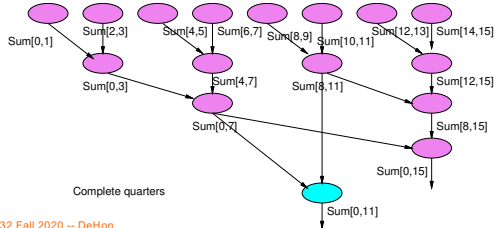


Penn ESE532 Fall 2020 -- DeHon

42

Prefix Tree

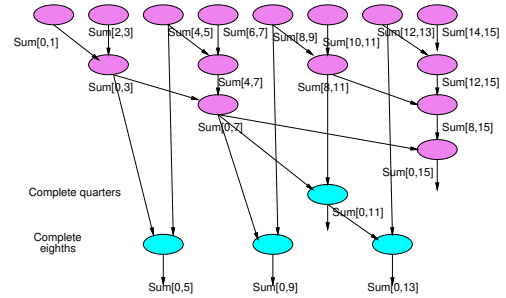
- Look at Symmetric stage (with respect to middle=Sum[0,N-1] stage) and combine to fill in



Penn ESE532 Fall 2020 -- DeHon

43

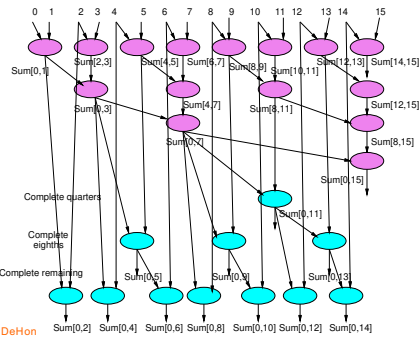
Prefix Tree



Penn ESE532 Fall 2020 -- DeHon

44

Prefix Tree

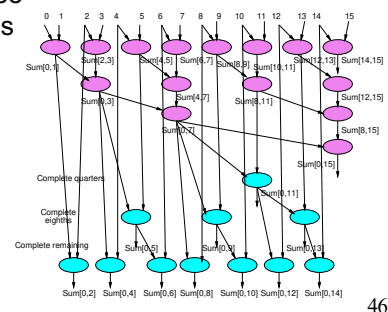


Penn ESE532 Fall 2020 -- DeHon

45

Prefix Tree

- Note: prefix-tree is same size as reduce tree

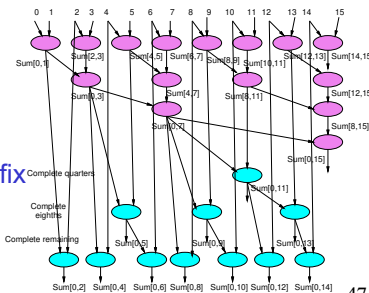


Penn ESE532 Fall 2020 -- DeHon

46

Parallel Prefix Area and Delay?

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



Penn ESE532 Fall 2020 -- DeHon

47

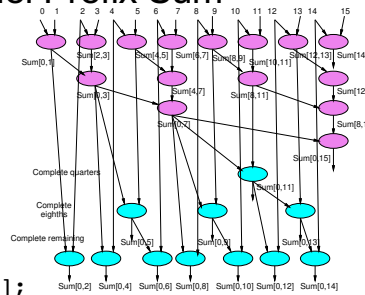
Parallel Prefix

- Important **Pattern**
- Applicable any time operation is *associative*
 - Or can be made associative
- Function Composition is always associative
 - (see Bonus at end)
- Logarithmic delay
- Linear area

Penn ESE532 Fall 2020 -- DeHon

48

Parallel Prefix Sum



```

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

BROADER APPLICATION

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)

Majority Associative?

- Carry=MAJ=majority
 - = $A \&\&B \parallel B \&\&C \parallel A \&\&C$
- Is Majority Associative ?
- Hint: What are each of following?
 - $MAJ(1,1,MAJ(1,1,MAJ(1,0,0)))$
 - $MAJ(MAJ(MAJ(1,1,1),1,1),0,0)$

Binary Addition

- Adding 2 W-bit numbers
 - What's the latency bound?
 - Area to achieve?
- boolean $a[i], b[i], s[i]$
- for ($i=0; i<W; i++$) {
 - $cn=(a[i]\&\&b[i]) \parallel (a[i]\&\&c) \parallel (b[i]\&\&c);$
 - $s[i]=a[i] \wedge b[i] \wedge c;$
 - $c=cn;$

Categorization

- To minimize confusion, will typically ask you to characterize:
 - Data parallel
 - Reduce
 - Sequential

Penn ESE532 Fall 2020 -- DeHon

55

Big Ideas:

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

Penn ESE532 Fall 2020 -- DeHon

56

Admin

- Feedback (including p2)
- ESE Talk Tuesday: Chris Batten
- No required reading for Wednesday
- P3 due Friday

Penn ESE532 Fall 2020 -- DeHon

57

Bonus

BINARY ADDITION

Penn ESE532 Fall 2020 -- DeHon

58

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;
}
```

Penn ESE532 Fall 2020 -- DeHon

59

## Associative

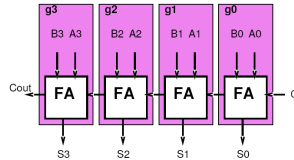
- Is the carry operation in addition an associative operation?
- Operation:
  - MAJ=majority =  $A \& B \ || \ B \& C \ || \ A \& C$

Penn ESE532 Fall 2020 -- DeHon

60

## 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  $a[i]$ ,  $b[i]$
  - $g=f(a,b)$

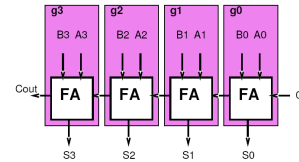


Penn ESE532 Fall 2020 -- DeHon

61

## Functions

- Carry=MAJ=majority  
 $= A\&B \parallel B\&C \parallel A\&C$
- What are the functions  $g(c[i-1])$ ?
  - $g(c)=\text{carry}(a=0,b=0,c)$
  - $g(c)=\text{carry}(a=1,b=0,c)$
  - $g(c)=\text{carry}(a=0,b=1,c)$
  - $g(c)=\text{carry}(a=1,b=1,c)$



Penn ESE532 Fall 2020 -- DeHon

62

## Functions

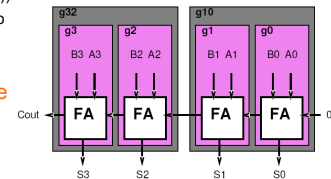
- What are the functions  $g(c[i-1])$ ?
  - $g(x)=1$  **Generate**
    - $a[i]=b[i]=1$
  - $g(x)=x$  **Propagate**
    - $a[i] \text{ xor } b[i]=1$
  - $g(x)=0$  **Squash**
    - $a[i]=b[i]=0$

Penn ESE532 Fall 2020 -- DeHon

63

## Combining

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



Penn ESE532 Fall 2020 -- DeHon

64

## Compose Rules (LSB MSB)

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

[work on board]

Penn ESE532 Fall 2020 -- DeHon

65

## Compose Rules (LSB MSB)

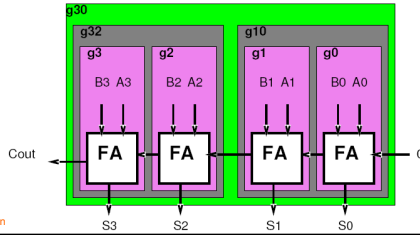
- GG = G
- GP = G
- GS = S
- PG = G
- PP = P
- PS = S
- SG = G
- SP = S
- SS = S

Penn ESE532 Fall 2020 -- DeHon

66

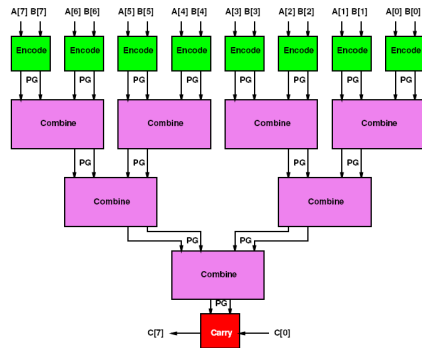
## Combining

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



Penn ESE532 Fall 2020 -- DeHon

## Associative Reduce Tree

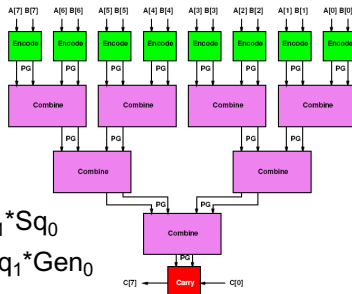


Penn ESE532 Fall 2020 -- DeHon

68

## Reduce Tree

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

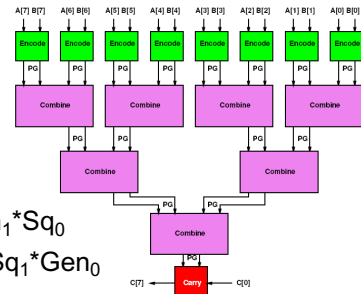


Penn ESE532 Fall 2020 -- DeHon

69

## 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 2020 -- DeHon

70

## Reduce Tree

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

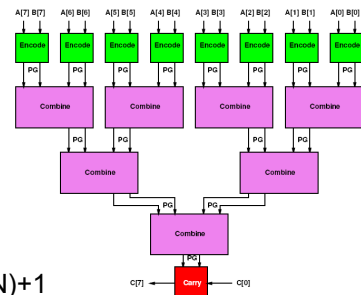
Penn ESE532 Fall 2020 -- DeHon

71

## Reduce Tree: Delay?

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

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



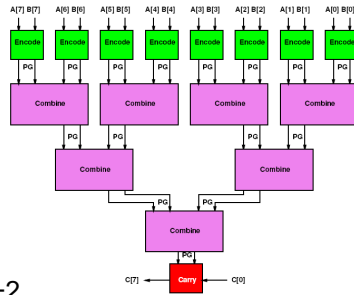
Penn ESE532 Fall 2020 -- DeHon

72

## Reduce Tree: Area?

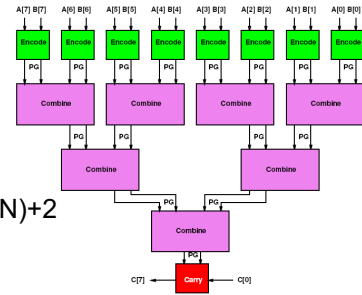
- $A(\text{Encode})=2$
- $A(\text{Combine})=4$
- $A(\text{Carry})=2$

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

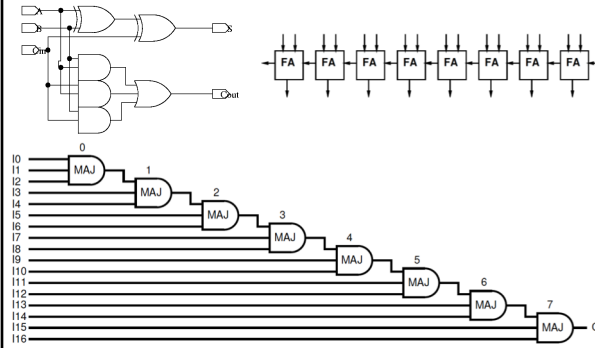


## Reduce Tree: Area & Delay

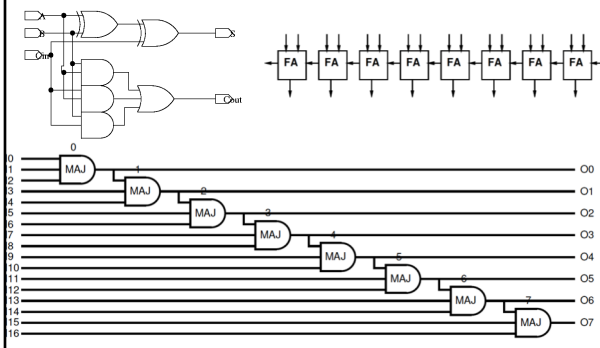
- $\text{Area}(N) = 6N - 2$
- $\text{Delay}(N) = 2\log_2(N) + 2$



## Compute Carry[N]



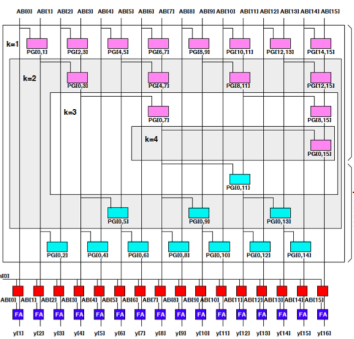
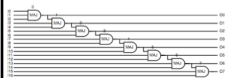
## Need Prefix



Need  $PG[i:0]$  for all  $i$

## Prefix Tree

- Bring in Carry and compute each intermediate carry-in



## Parallel Prefix Area and Delay?

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

