

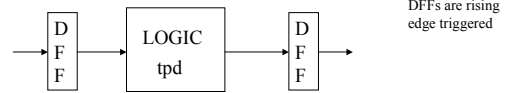
## Speeding up the Clock

- The register-to-register delay is usually the delay path that sets the maximum clock rate
- From a design point of view, can only affect the combinational logic between the registers
  - Need to shorten the maximum combinational delay path
  - Setup/Hold time of registers are fixed
- Can shorten the delay by placing a register in the combinational logic to break longest delay path
  - This technique is called *pipelining*
  - Adds *latency* to the output (the number of clocks between an input value and its corresponding output result)

BR 1/99

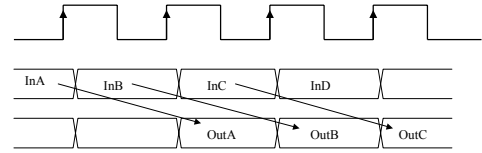
1

## Registered Datapath



$$\text{Clk Freq} = 1 / (\text{Telk}2q + \text{Tpd} + \text{Tsu})$$

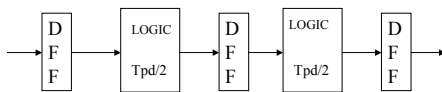
Latency = 1 clk



BR 1/99

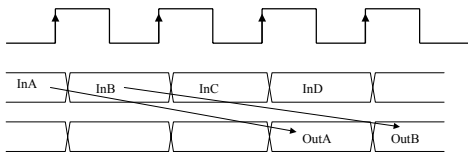
2

## Add a pipeline stage



$$\text{Clk Freq} = 1 / (\text{Telk}2q + \text{Tpd}/2 + \text{Tsu})$$

Latency = 2 clks



BR 1/99

3

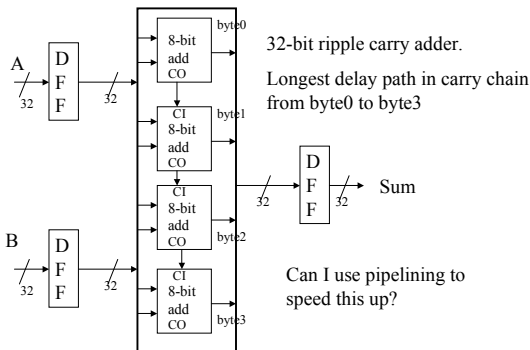
## Definitions

- Initiation Rate** - Rate at which new input values are accepted after first output value
  - Rate at which new computations are initiated
  - Maximum initiation rate = 1
- Latency**
  - Number of clock cycles between input value and output value
  - Adding pipeline stages always increases latency
- At some point, adding more pipeline stages does not increase clock frequency because  $\text{Telk}2q$  and  $\text{Tsu}$  dominate delay.

BR 1/99

4

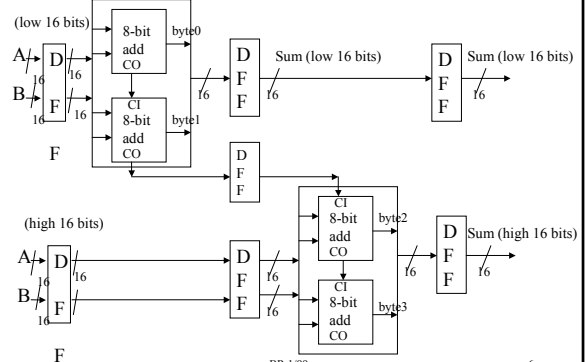
## Pipeline Example



BR 1/99

5

## Insert pipeline stage between byte1 and byte2.



BR 1/99

6

## Comments on Pipeline Example

- Note that the pipeline stage broke the carry chain into two equal paths
  - Each pipeline stage should have approximately the same combinational delay
  - Clock speed will be set by the delay of the slowest pipeline stage
- If I inserted 2 pipeline stages, I would need to break the carry chain delay into equal thirds
- Could insert a pipeline stage between each BIT in order to get maximum clock speed
  - Called 'bit pipelining'

BR 1/99

7

## Latency Tolerance

- Latency tolerance is dependent upon each application
- Frequent flushing of a pipeline (discarding partial results within the pipeline and restarting the pipeline with a new value) wastes time and makes an application latency **intolerant**.
- Flushing of a pipeline introduces clock cycles in which the results coming out of the pipeline are ignored -- these are wasted clock cycles.
- High Latency tolerance means that you can have many pipeline stages, whatever the number you need to meet the clock rate specification.

BR 1/99

8

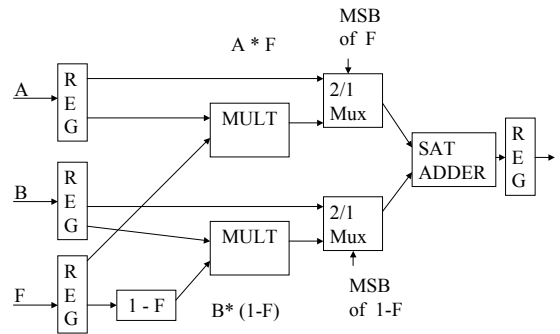
## Two Applications

- Graphics hardware for processing pixels is extremely latency tolerant - not unusual to find pipelines that have 10's of stages.
  - Graphics pipelines are never flushed
  - High clock rate is EXTREMELY important because of large number of pixels (> 1 Million) that have to be supplied every screen, at > 30 updates per second
- Microprocessor instruction pipelines are not very latency tolerant - most CPU pipelines are only about 5-10 stages.
  - Branch instructions can cause pipeline to be flushed. By the time you determine direction of branch, may have started processing instructions that should not be in the pipeline. These are flushed and the pipeline restarted.

BR 1/99

9

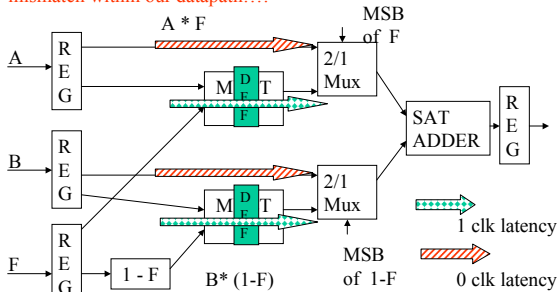
BLEND Datapath without Pipelining. MULT is combinational.



BR 1/99

10

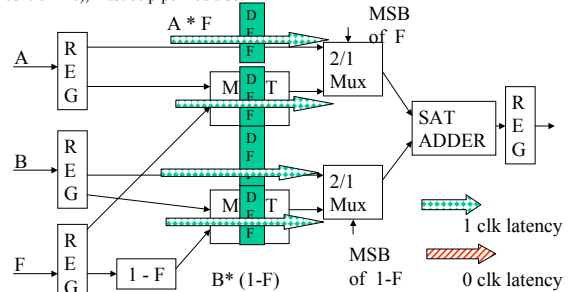
LPM\_Mult can be pipelined via LPM\_PIPELINE parameter. Add one pipeline stage to LPM\_MULT. DFFs inserted automatically in LPM\_MULT. **We now have a LATENCY mismatch within our datapath!!!!**



BR 1/99

11

Correct the latency mismatch by adding DFFs in other path as well. May have to break delay paths in other places or add additional pipeline stages to LPM\_mult to meet clock frequency target. Control lines (like MUX control line), must be pipelined also.



BR 1/99

12