Parallelism

- Two types of parallelism:
 - Spatial parallelism
 - duplicate hardware performs multiple tasks at once
 - Temporal parallelism
 - task is broken into multiple stages
 - also called pipelining
 - for example, an assembly line



Parallelism Definitions

- **Token:** Group of inputs processed to produce group of outputs
- Latency: Time for one token to pass from start to end
- **Throughput:** Number of tokens produced per unit time

Parallelism increases throughput



Parallelism Example

- Ben Bitdiddle bakes cookies to celebrate traffic light controller installation
- 5 minutes to roll cookies
- 15 minutes to bake
- What is the latency and throughput without parallelism?



Parallelism Example

- Ben Bitdiddle bakes cookies to celebrate traffic light controller installation
- 5 minutes to roll cookies
- 15 minutes to bake
- What is the latency and throughput without parallelism?

Latency = 5 + 15 = 20 minutes = 1/3 hour Throughput = 1 tray/ 1/3 hour = 3 trays/hour



Parallelism Example

- What is the latency and throughput if Ben uses parallelism?
 - Spatial parallelism: Ben asks Allysa P. Hacker to help, using her own oven
 - Temporal parallelism:
 - two stages: rolling and baking
 - He uses two trays
 - While first batch is baking, he rolls the second batch, etc.



Spatial Parallelism





Spatial Parallelism



Latency = 5 + 15 = 20 minutes = 1/3 hour Throughput = 2 trays/ 1/3 hour = 6 trays/hour



Temporal Parallelism



Latency = ? Throughput = ?



© Digital Design and Computer Architecture, 2nd Edition, 2012

Chapter 3 <89>

Temporal Parallelism



Latency = 5 + 15 = 20 minutes = 1/3 hour Throughput = 1 trays/ 1/4 hour = 4 trays/hour

Using both techniques, the throughput would be 8 trays/hour

