



Real-Time Systems

Lecture Topic - Accounting for Execution Efficiency

Dr. Sam Siewert

Electrical, Computer and Energy Engineering

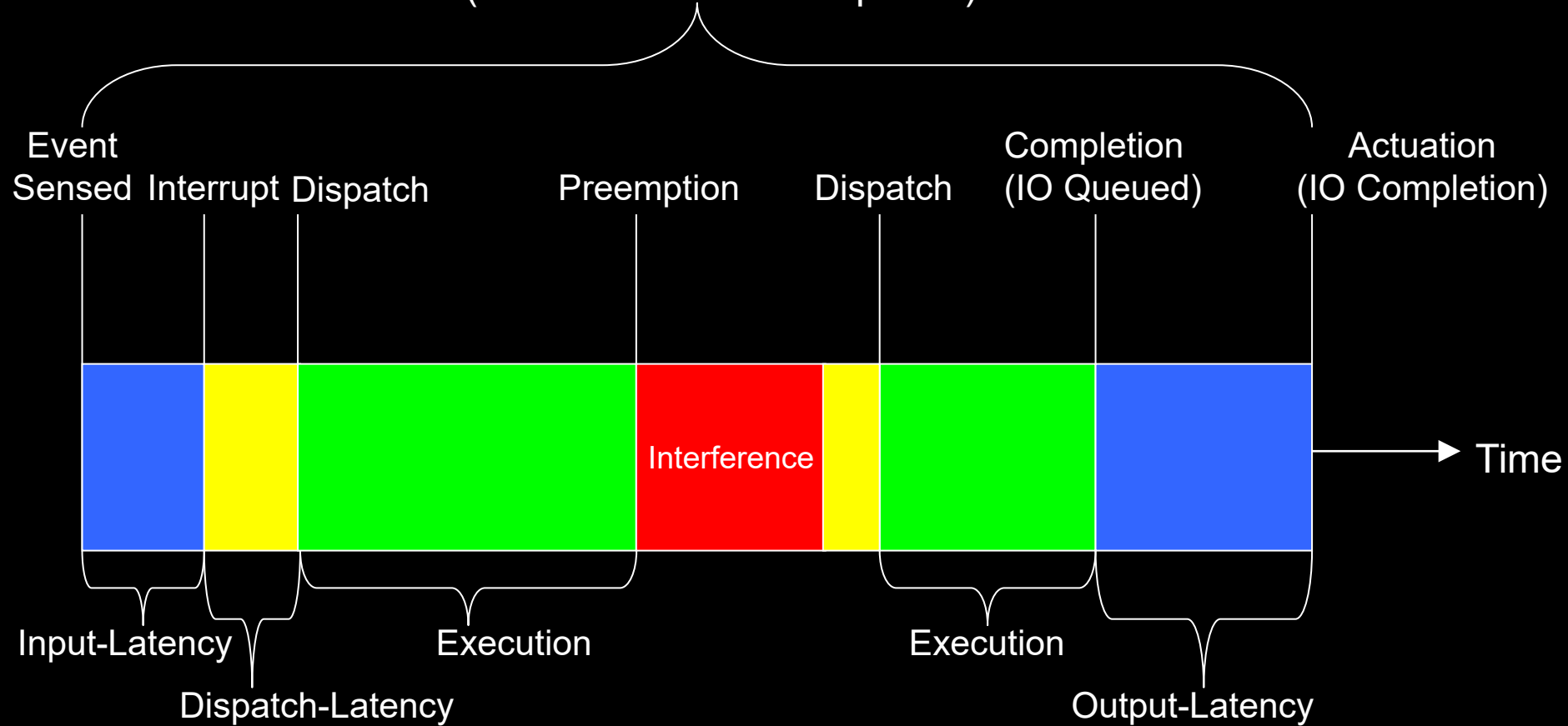
Embedded Systems Engineering Program

A Service Release and Response

- C_i WCET
- Input/Output Latency
- Interference Time

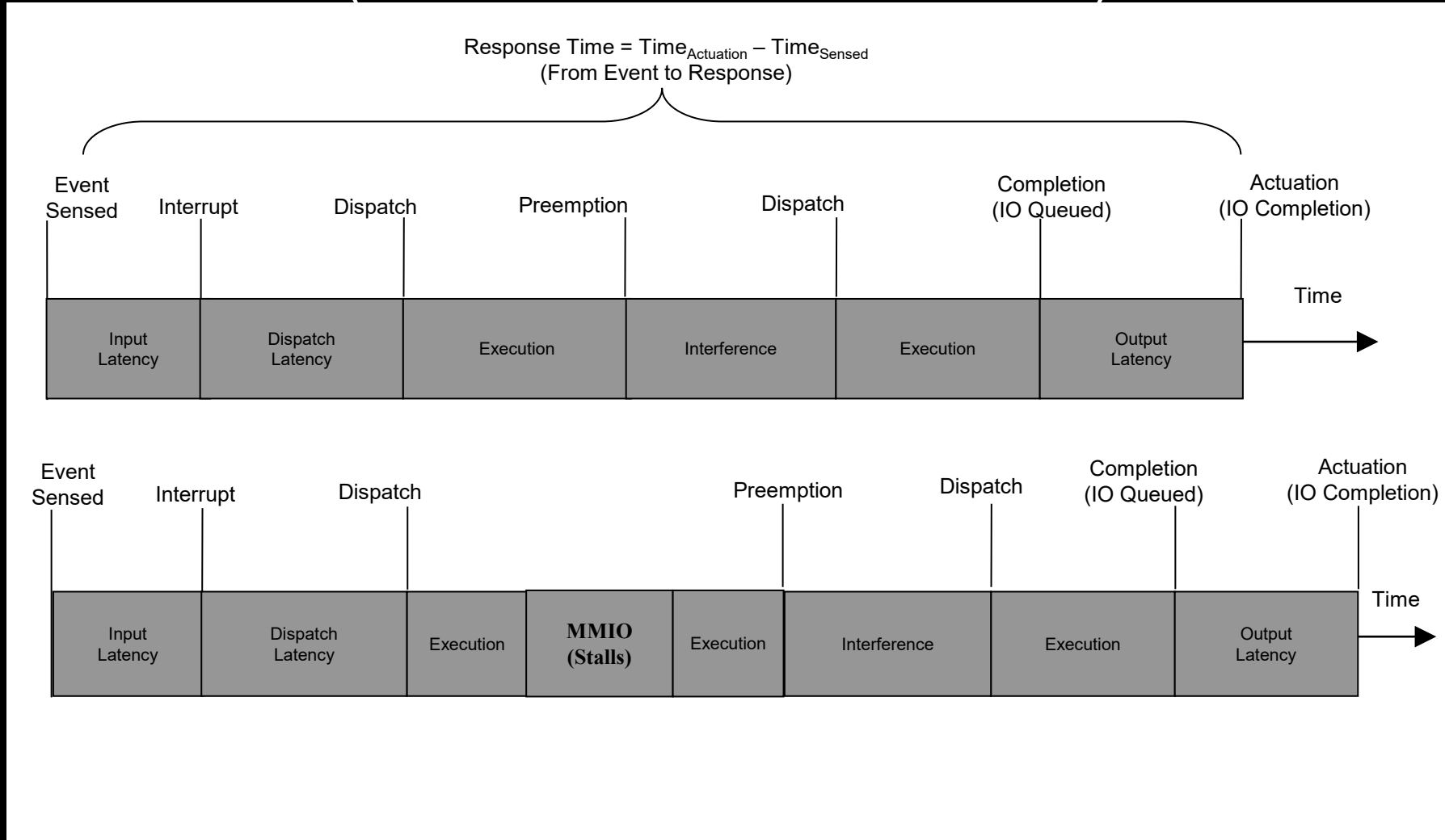
$$\text{Response Time} = \text{Time}_{\text{Actuation}} - \text{Time}_{\text{Sensed}}$$

(From Release to Response)



Service Response Timeline - WCET

(With Intermediate I/O – WCET)

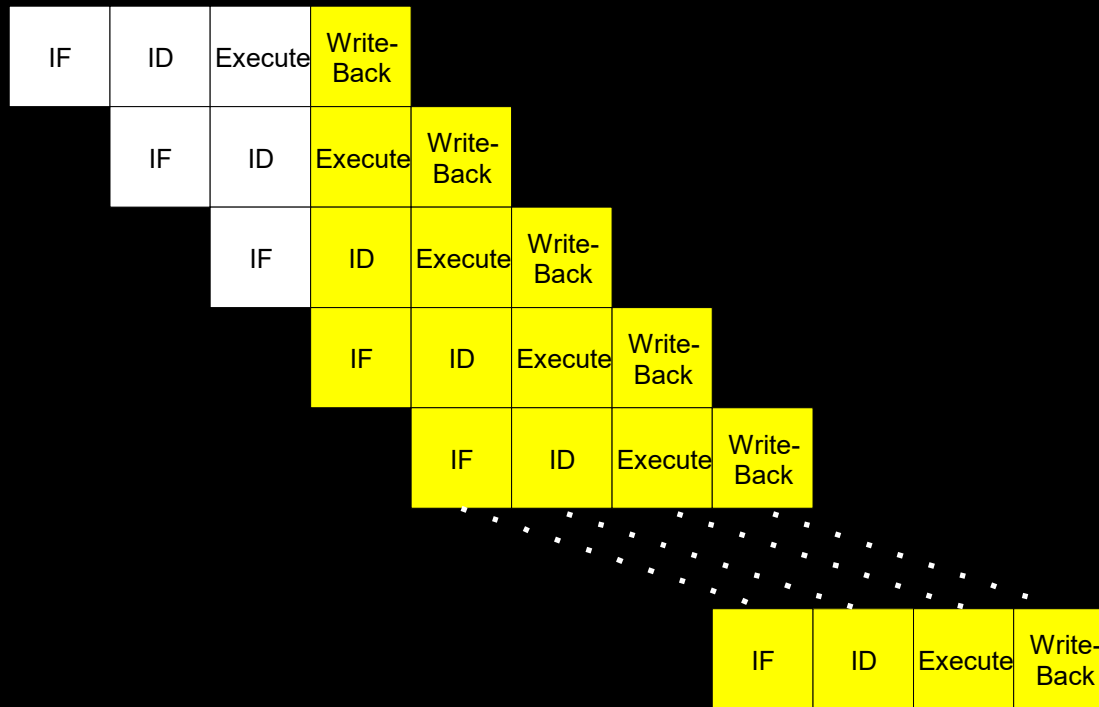


Services and Execution Efficiency

- CPU Pipeline Stalls (Hazards)
 - Cache Miss/Hit
 - Off-Chip MMIO
 - Branch Misprediction
 - Data Dependencies
 - Register Pressure, MMR (Memory Mapped Registers)
- During Service Execution, Memory Mapped Register I/O
- External Memory I/O
- Cache Loads and Write-Backs
- In Memory Input/Output from Service to Service

Pipelined Architecture Review

- Recall that Pipeline Yields CPI of 1 or Less
- Instruction Completed Each CPU Clock
- Unless Pipeline Stalls!



Service Execution

$$WCET = [(CPI_{best-case} \times Longest_Path_Inst_Count) + Stall_Cycles] \times Clk_Period$$

- Efficiency of Execution
- Memory Access for Inter-Service Communication
- Bounded Intermediate I/O
- Ideally Lock Data and Code into Cache or Use Tightly Coupled Memory [Scratch Pad, Zero Wait State]

Service Efficiency

■ Path Length for a Release

- Instruction Count
- Longest Path Given Algorithm
 - Branches
 - Loop Iterations
 - Data Driven?

■ Path Execution

- Number of Cycles to Complete Path
- Clocks Per Instruction
- Number of Stall Cycles for Pipeline
 - Data Dependencies (Intermediate IO)
 - Cache Misses (Intermediate Memory Access)
 - Branch Mis-predictions (Small Penalty)

Hiding Intermediate IO Latency

(Overlapping CPU and Bus I/O)

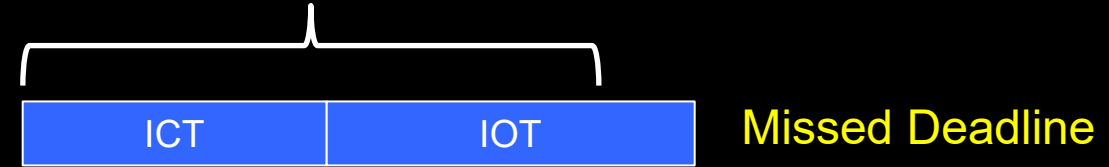
- **ICT = Instruction Count Time**
 - $(CPI=1) \times Inst_Cnt \times CPU_Clk_Period$
 - Time to execute a block of instructions with no stalls, hence $CPI=1$
 - $CPU_Clk_Cycles \times CPU_Clock_Period$
- **IOT = Bus Interface IO Time**
 - $Bus\ IO\ Cycles \times Bus\ Clock\ Period$
 - $Bus\ Clock\ Period / CPU_Clk_Period = Wait\ Clocks$
- **OR = Overlap Required**
 - Percentage of CPU cycles that must be concurrent with I/O cycles
- **NOA = Non-Overlap Allowable for S_i to meet D_i**
 - Percentage of CPU cycles that can be in addition to IO cycle time without missing service deadline
- **D_i = Deadline for Service S_i relative to release**
 - interrupt or system call initiates S_i request and S_i execution
- **CPI = Clocks Per Instruction for a block of instructions**
 - IPC is Instructions Per Clock, Also Used, Just the Inverse (Superscalar, Pipelined)

Processing and I/O Time Overlap

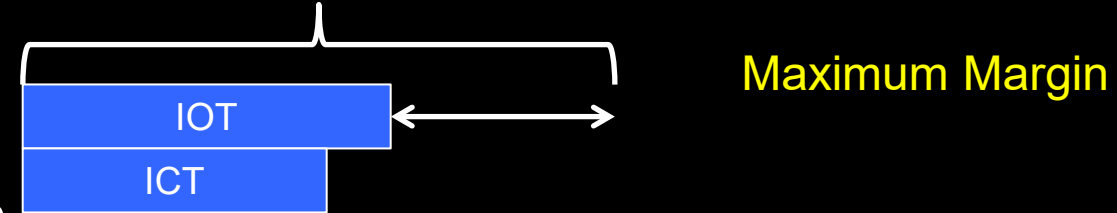
- If $D_i < IOT$, S_i is **IO-Bound**
Hardware Re-design REQUIRED
- If $D_i < ICT$, S_i is **CPU-Bound**
- if $D_i < (IOT + ICT)$ where ($D_i < IOT$ AND $D_i < ICT$), deadline **D_i can't be met regardless of overlap, both IO-Bound and CPU-Bound**
- $D_i \geq (IOT + ICT)$ **requires no overlap** of IOT with ICT
No code I/O optimization required
- if $D_i < (IOT + ICT)$ where ($D_i \geq IOT$ and $D_i \geq ICT$), **overlap of IOT with ICT is required**
Code I/O optimization required
- $CPI_{ICT} = 1$ by definition for ICT alone (no stalls, best case efficiency)
- $CPI_{actual} = f_{efficiency} \times CPI_{ICT}$

Service Execution Efficiency – Waiting on Bus and Memory I/O

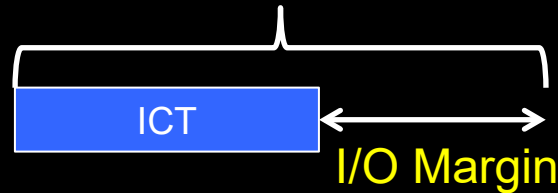
- $CPI_{\text{worst-case}} = (ICT + IOT) / ICT$



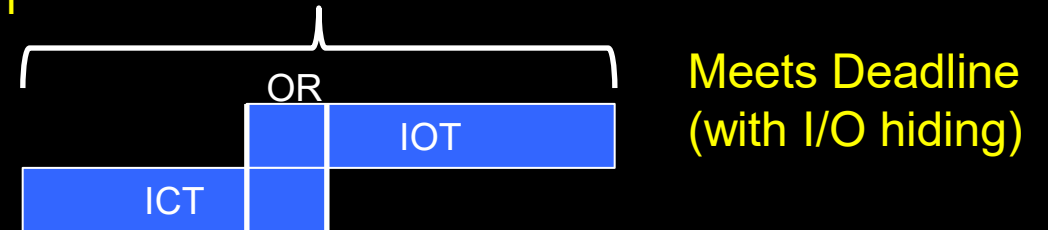
- $CPI_{\text{best-case}} = (\max(ICT, IOT)) / ICT$



- $CPI_{\text{required}} = D_i / ICT$



- $OR = 1 - [(D_i - IOT) / ICT]$



- $CPI_{\text{required}} = [ICT(1-OR) + IOT] / ICT$

- $NOA = (D_i - IOT) / ICT; OR + NOA = 1 \text{ (by definition)}$

CPI Units

■ $\text{CPI}_{\text{required}} = D_i / \text{ICT} = \text{Time}/\text{Time}?$

- $D_i = \text{Clocks} \times \text{Clock_T} = \text{Time} [\text{N} \times \text{X nanoseconds}]$
- $\text{ICT} = \text{CPI}=1 \times \text{Num_instructions} \times \text{Clock_T} = \text{Time} [\text{N} \times \text{X nanoseconds}]$

■ $D_i / \text{ICT} = (\text{Clocks} \times \text{Clock_T}) / (1 \times \text{Num_instructions} \times \text{Clock_T})$

- $\text{Clocks} / \text{Num_instructions} = \text{CPI}_{\text{required}}$ over deadline D_i

■ $\text{IPC} = 1/\text{CPI}$ (IPC used for pipelined, superscalar)

Processing Latency Alone

- Write Code with Memory Resident Frames
 - Load Frames in Advance
 - Process In-Memory Frames Over and Over
 - Do No I/O During Processing
 - Provides Baseline Measurement of Processing Latency per Frame Alone
 - Provides Method of Optimizing Processing Without IO Latency

