# Dynamic Tuning of Feature Set in Highly Variant Interactive Applications

**Tushar Kumar**, Romain Cledat,

and Santosh Pande

**Georgia Tech** | College of Computing

*EMSOFT 2010*

# Motivation

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications " *lack analyzable semantics*

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications " *lack analyzable semantics*
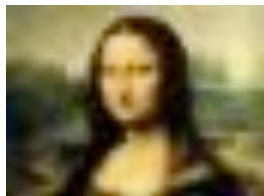
**30 fps**

Per-Frame Time

Slow Processor

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications " *lack analyzable semantics*

**30 fps**

Per-Frame Time

Slow Processor

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications " *lack analyzable semantics*
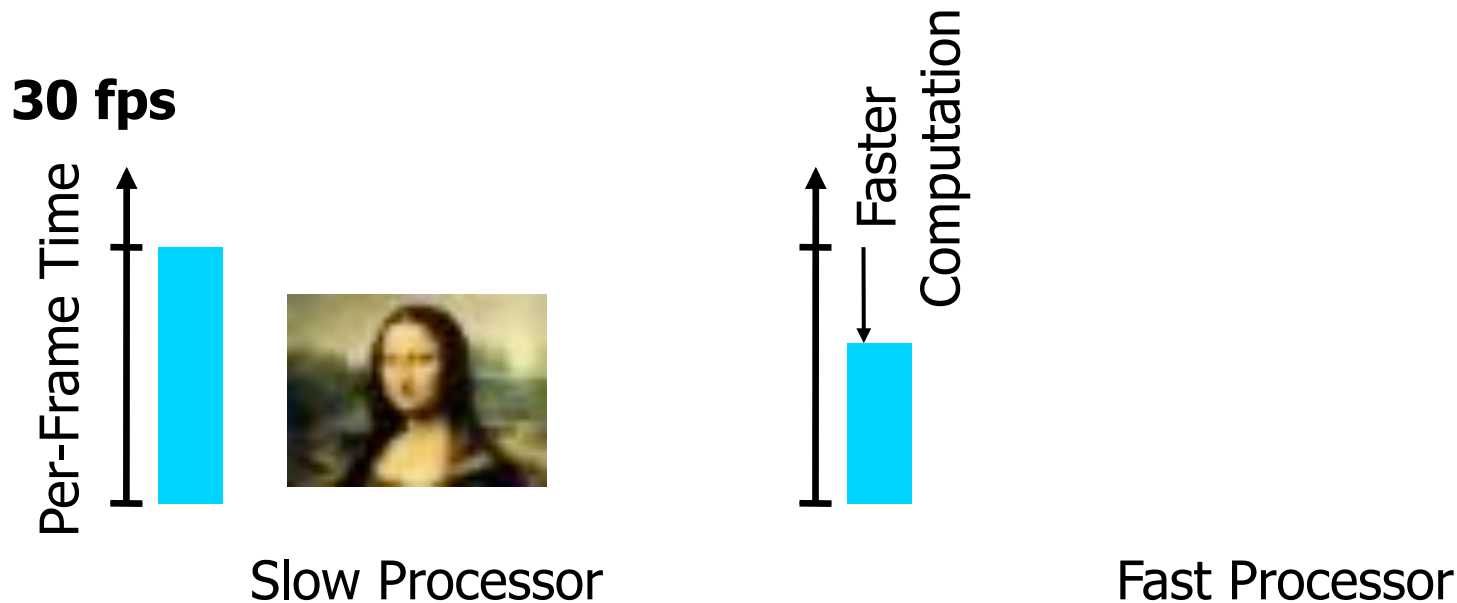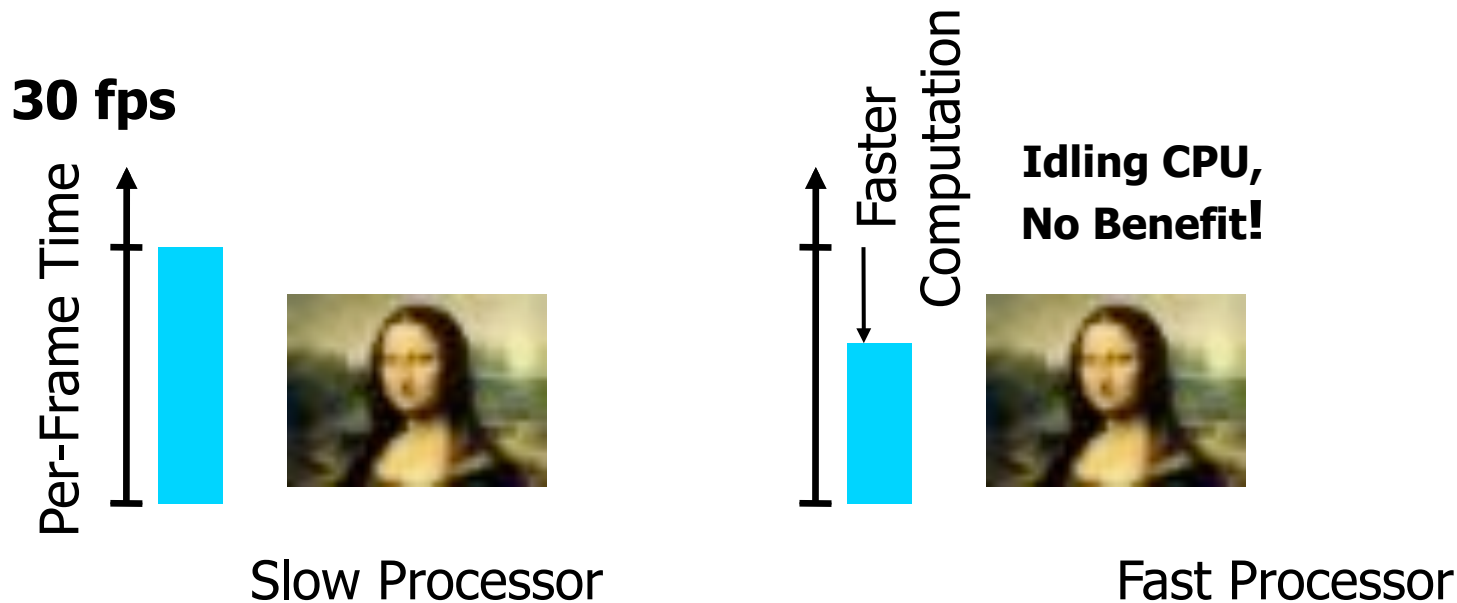


**30 fps**

Per-Frame Time

Faster Computation

Slow Processor

Fast Processor

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications '' *lack analyzable semantics*
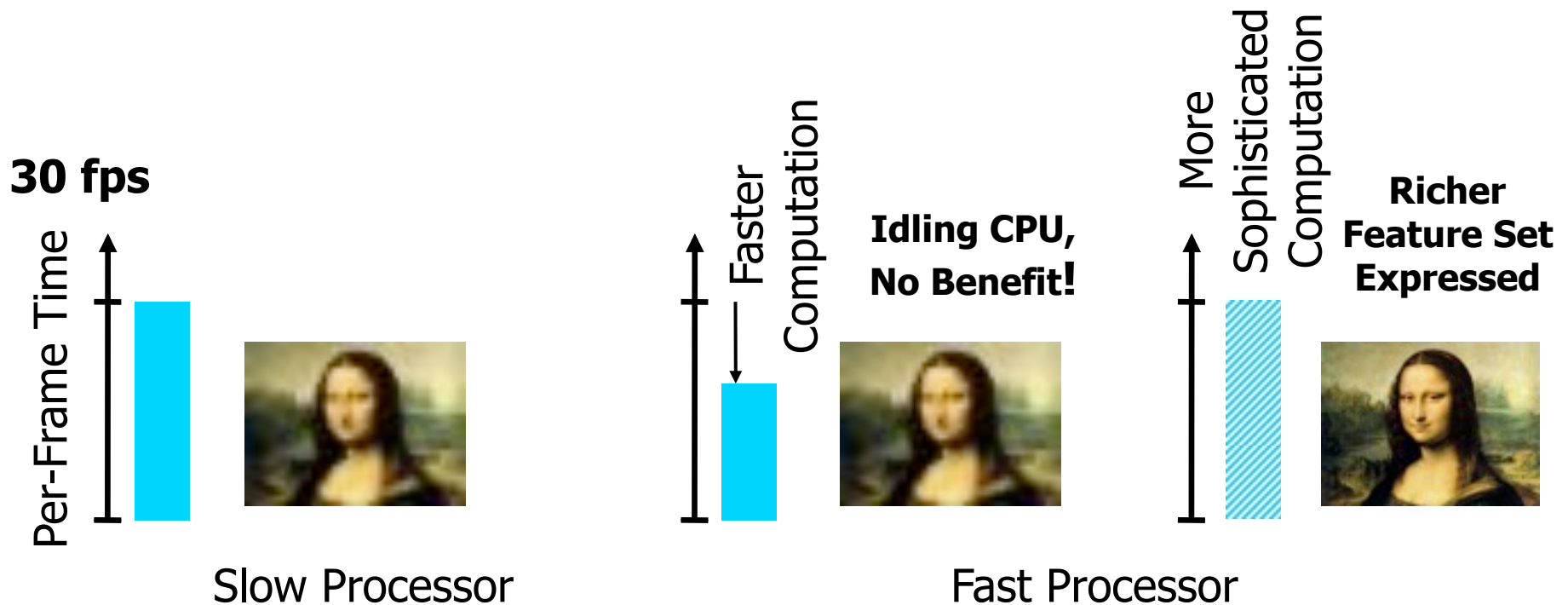


**30 fps**

Per-Frame Time

Slow Processor

Faster Computation

**Idling CPU, No Benefit!**

Fast Processor

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications " *lack analyzable semantics*
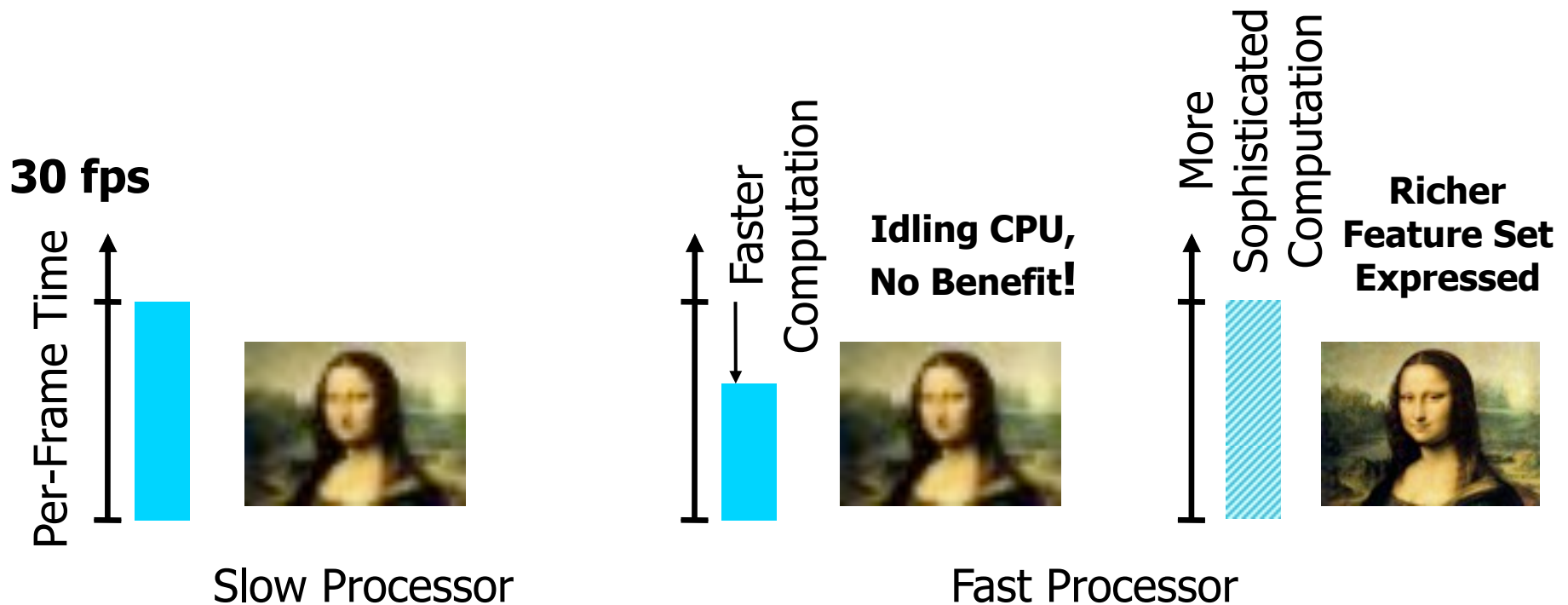


**30 fps**

Per-Frame Time

Faster Computation

**Idling CPU, No Benefit!**

More Sophisticated Computation

**Richer Feature Set Expressed**

Slow Processor

Fast Processor

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications " *lack analyzable semantics*
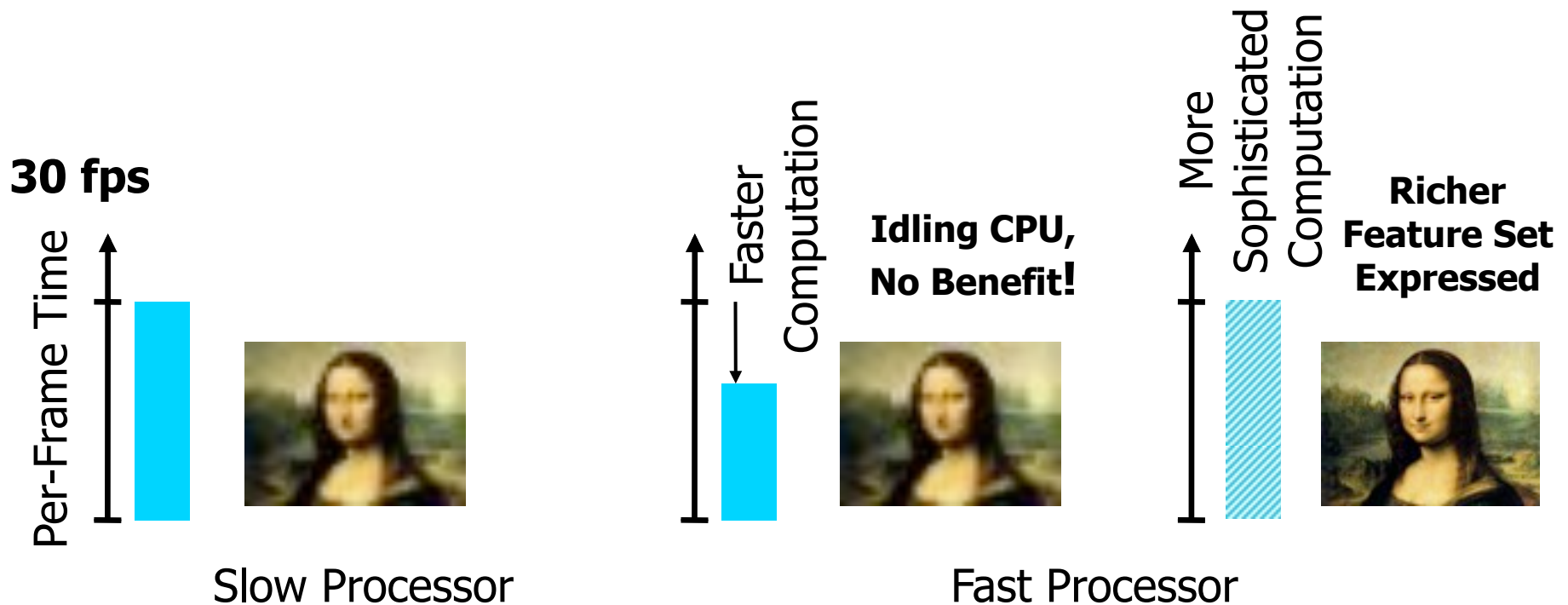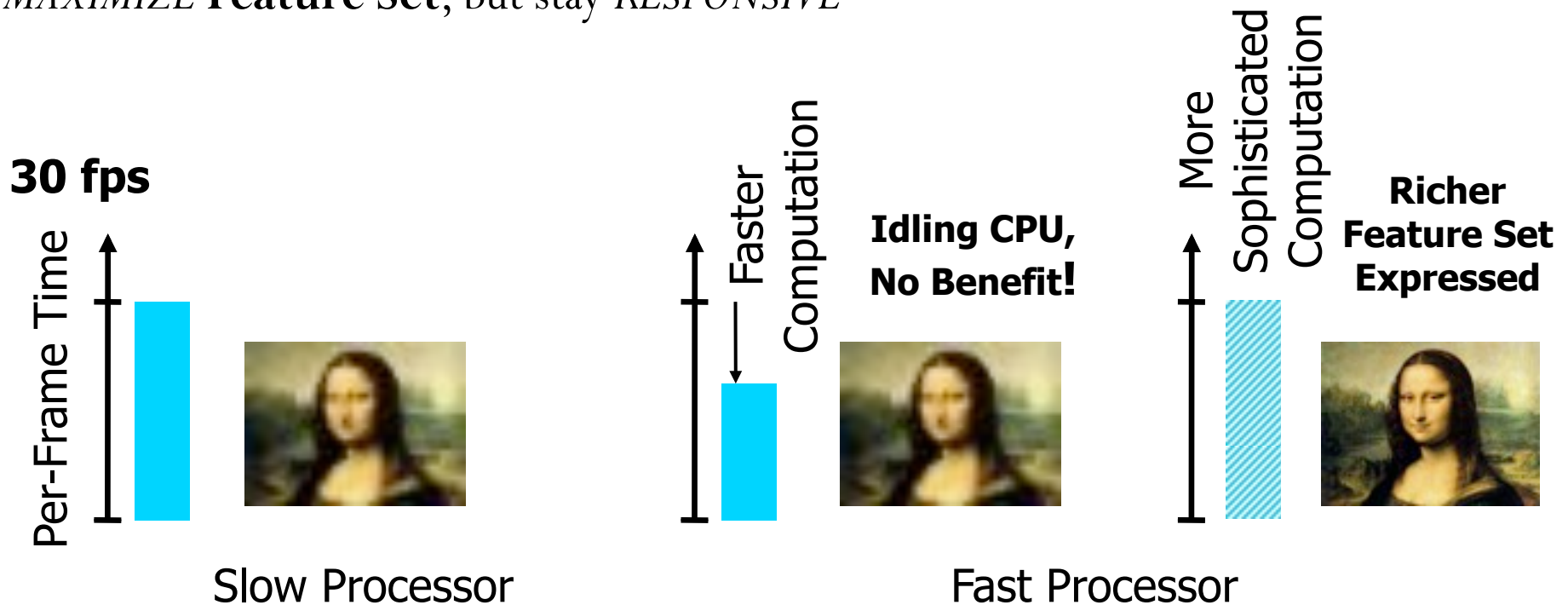
- Design goal



**30 fps**

Per-Frame Time

Slow Processor

Faster Computation

Idling CPU, No Benefit!

Fast Processor

More Sophisticated Computation

**Richer Feature Set Expressed**

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications " *lack analyzable semantics*

- Design goal
  - Provide the **richest, most engrossing experience** possible to the **interactive user**



**30 fps**

Per-Frame Time

Faster Computation

**Idling CPU, No Benefit!**

More Sophisticated Computation

**Richer Feature Set Expressed**

Slow Processor

Fast Processor

# Motivation

- An entire class of **frame-oriented**, **interactive applications** currently lack a systematic optimization methodology
  - Gaming, Multimedia, Interactive Visualization
  - Conventionally developed C/C++/Java applications '' *lack analyzable semantics*

- Design goal
  - Provide the **richest, most engrossing experience** possible to the **interactive user**
  - *MAXIMIZE* **Feature Set**, but stay *RESPONSIVE*



**30 fps**

Per-Frame Time

Slow Processor

Faster Computation

Idling CPU, No Benefit!

More Sophisticated Computation

Richer Feature Set Expressed

Fast Processor

# Representative Applications

# Representative Applications

- MPEG2 Encoder

# Representative Applications

- ## MPEG2 Encoder
  - Motion Estimation algorithm dominates per-frame time
  - ***Search-Window-Size*** parameter dramatically ***scales***:
    - Per-frame encoding time

# Representative Applications

- ## MPEG2 Encoder
  - Motion Estimation algorithm dominates per-frame time
  - *Search-Window-Size* parameter dramatically *scales*:
    - Per-frame encoding time
    - Achieved compression of video

# Representative Applications

- ## MPEG2 Encoder
  - Motion Estimation algorithm dominates per-frame time
  - ***Search-Window-Size*** parameter dramatically ***scales***:
    - Per-frame encoding time
    - Achieved compression of video

- ## Torque Game Engine (www.torquepowered.com)
  - AI algorithms dominate per-frame time, when there are large number of simulated enemies (bots)
  - Multiple AI parameters ***scale***:

# Representative Applications

- ## MPEG2 Encoder
  - Motion Estimation algorithm dominates per-frame time
  - *Search-Window-Size* parameter dramatically *scales*:
    - Per-frame encoding time
    - Achieved compression of video

- ## Torque Game Engine (www.torquepowered.com)
  - AI algorithms dominate per-frame time, when there are large number of simulated enemies (bots)
  - Multiple AI parameters *scale*:
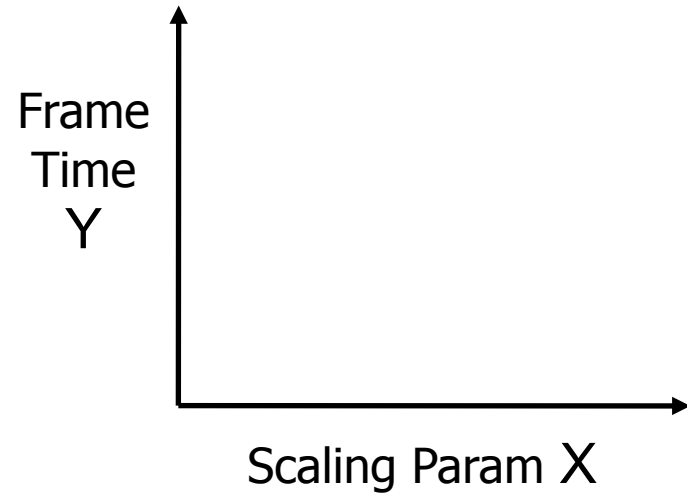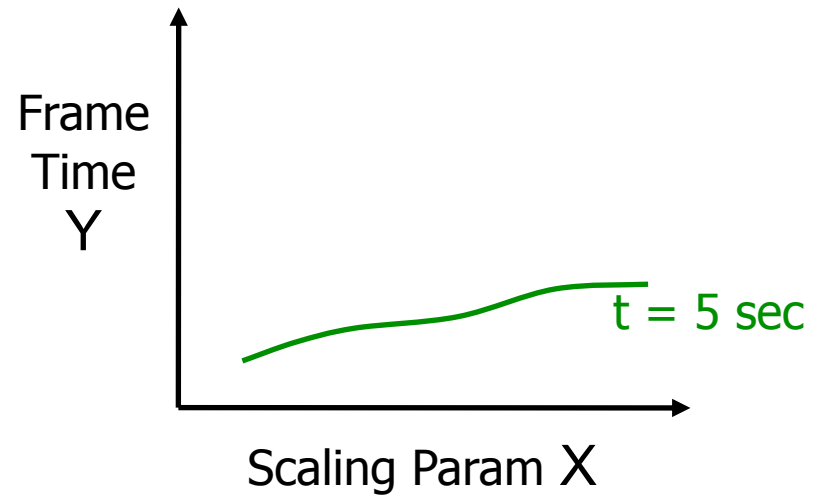    - Frame-rate
    - Intelligence exhibited by bots

# Representative Applications

- ## MPEG2 Encoder
  - Motion Estimation algorithm dominates per-frame time
  - *Search-Window-Size* parameter dramatically *scales*:
    - Per-frame encoding time
    - Achieved compression of video

- ## Torque Game Engine (www.torquepowered.com)
  - AI algorithms dominate per-frame time, when there are large number of simulated enemies (bots)
  - Multiple AI parameters *scale*:
    - Frame-rate
    - Intelligence exhibited by bots

- Scalable Algorithms are common and dominant

# Representative Applications

- # MPEG2 Encoder
  - Motion Estimation algorithm dominates per-frame time
  - *Search-Window-Size* parameter dramatically *scales*:
    - Per-frame encoding time
    - Achieved compression of video

- # Torque Game Engine (www.torquepowered.com)
  - AI algorithms dominate per-frame time, when there are large number of simulated enemies (bots)
  - Multiple AI parameters *scale*:
    - Frame-rate
    - Intelligence exhibited by bots

- Scalable Algorithms are common and dominant
  - Let X denote the *scaling parameter* (application-specific choice)

# Challenges in Tuning Feature Set
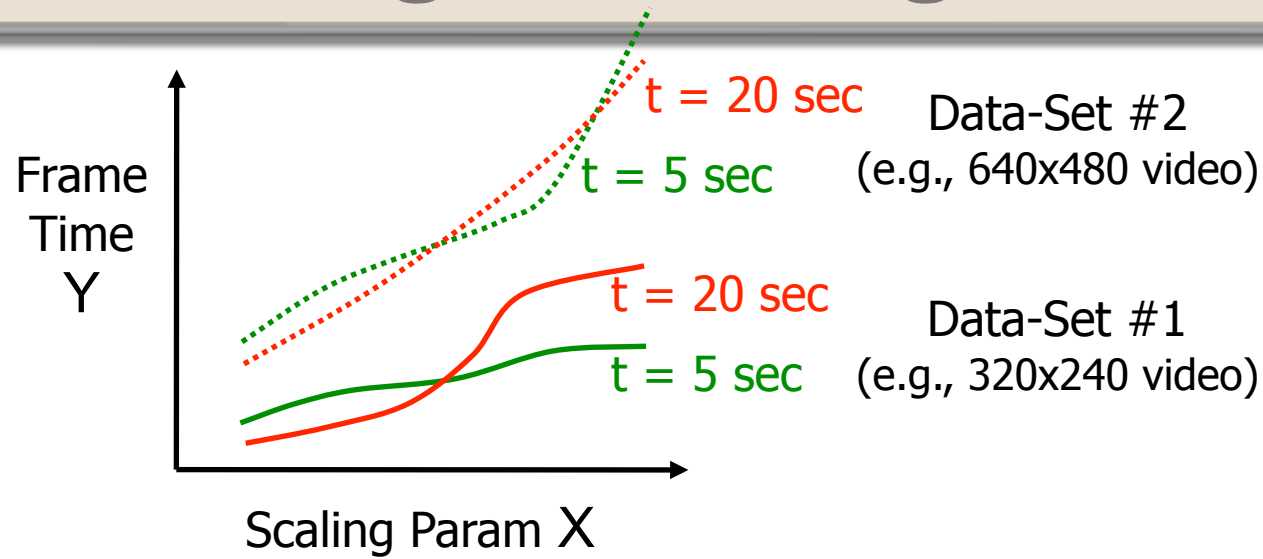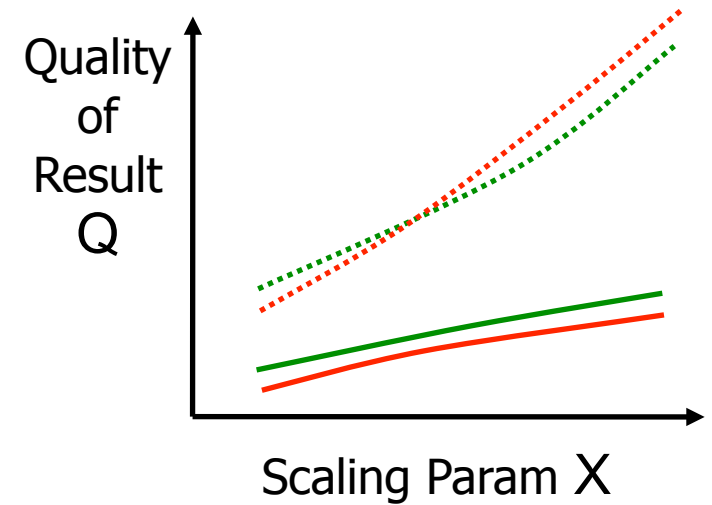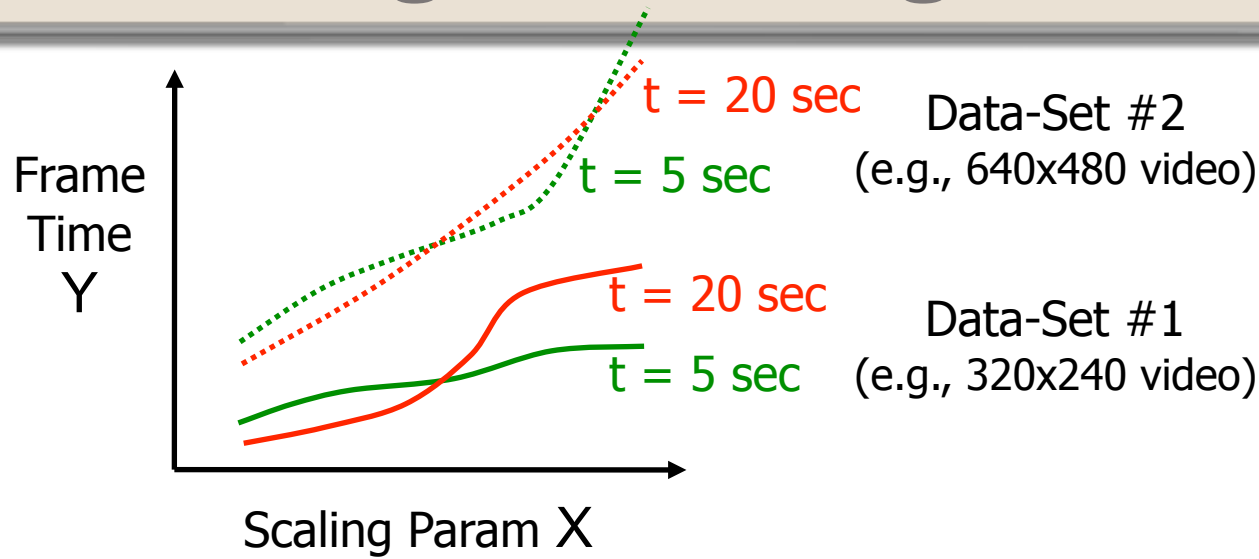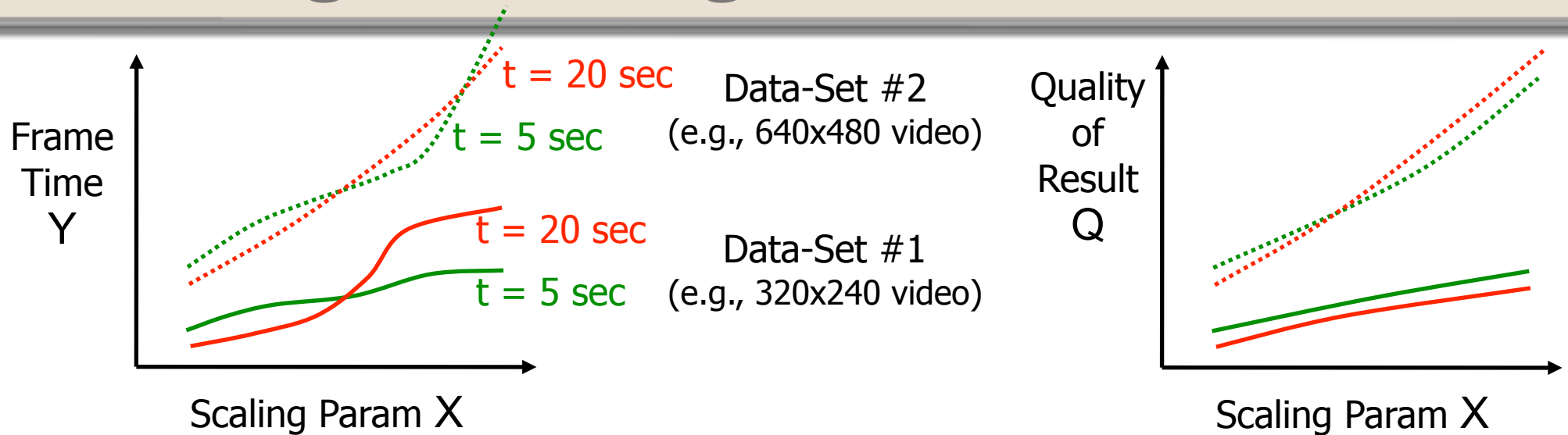
# Challenges in Tuning Feature Set

# Challenges in Tuning Feature Set

# Challenges in Tuning Feature Set



Frame Time Y

t = 20 sec

t = 5 sec

Data-Set #1
(e.g., 320x240 video)

Scaling Param X

# Challenges in Tuning Feature Set



Frame
Time
Y

t = 20 sec
Data-Set #2
t = 5 sec    (e.g., 640x480 video)

t = 20 sec
Data-Set #1
t = 5 sec    (e.g., 320x240 video)

Scaling Param X

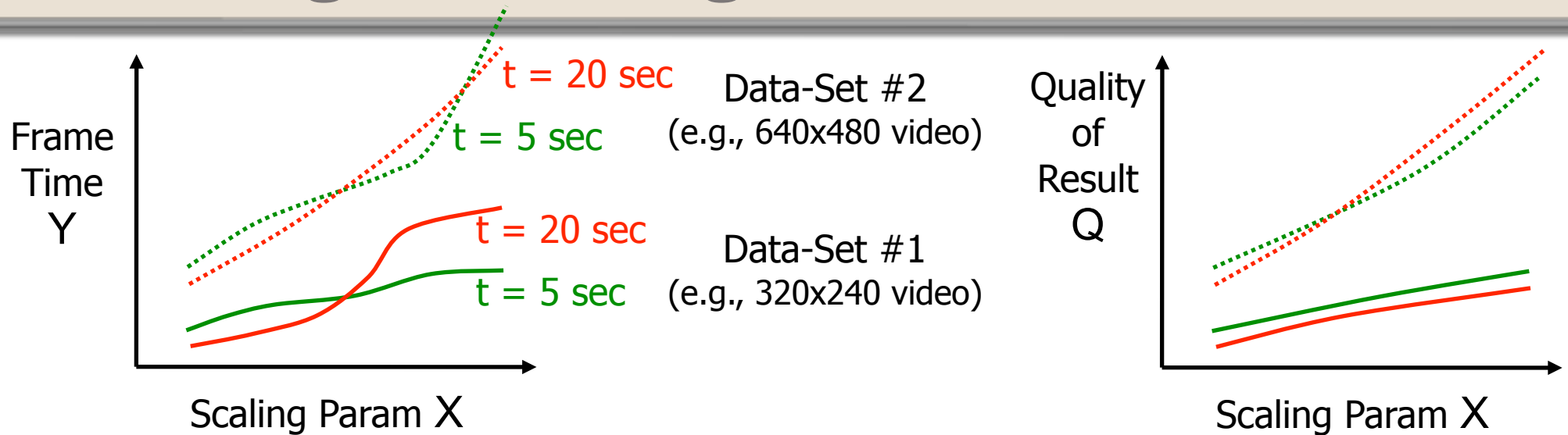# Challenges in Tuning Feature Set

# Challenges in Tuning Feature Set



X-Y, X-Q Response Characteristics
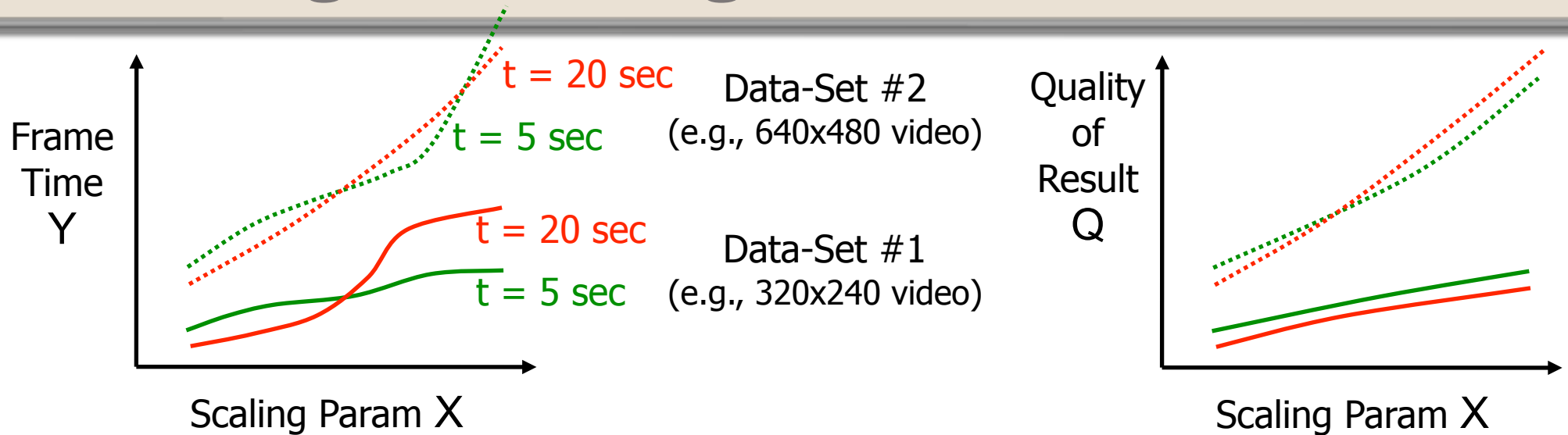- Highly time-varying, data-set dependent ” *No fixed relationships*

# Challenges in Tuning Feature Set



X-Y, X-Q Response Characteristics

- Highly time-varying, data-set dependent " *No fixed relationships*
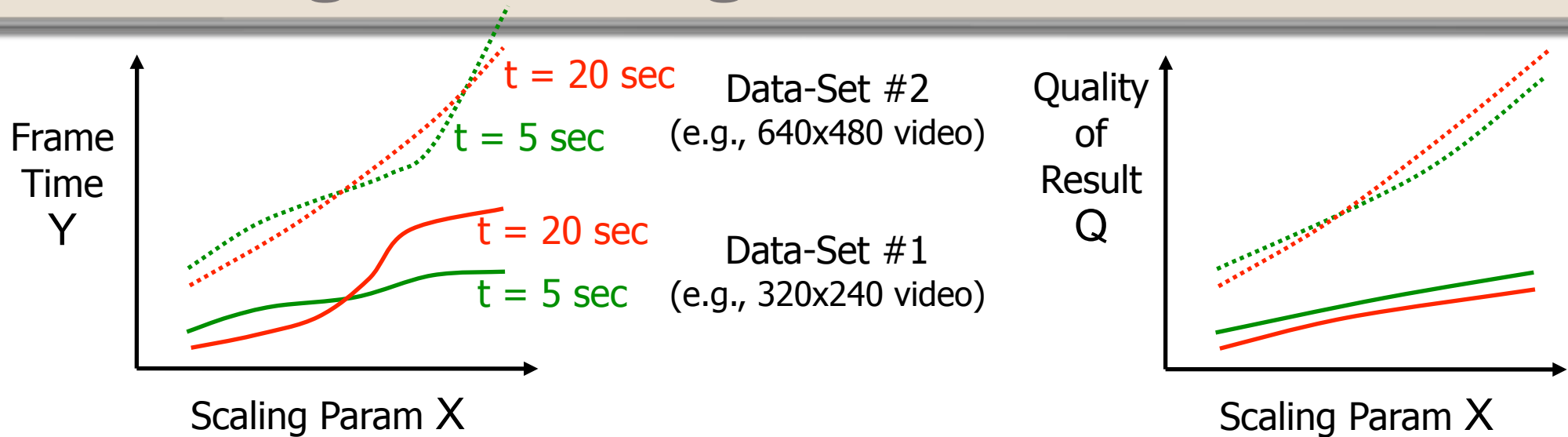- For practical purposes: UNKNOWN, entirely EMERGENT behavior

# Challenges in Tuning Feature Set



X-Y, X-Q Response Characteristics

- Highly time-varying, data-set dependent ” *No fixed relationships*
- For practical purposes: UNKNOWN, entirely EMERGENT behavior
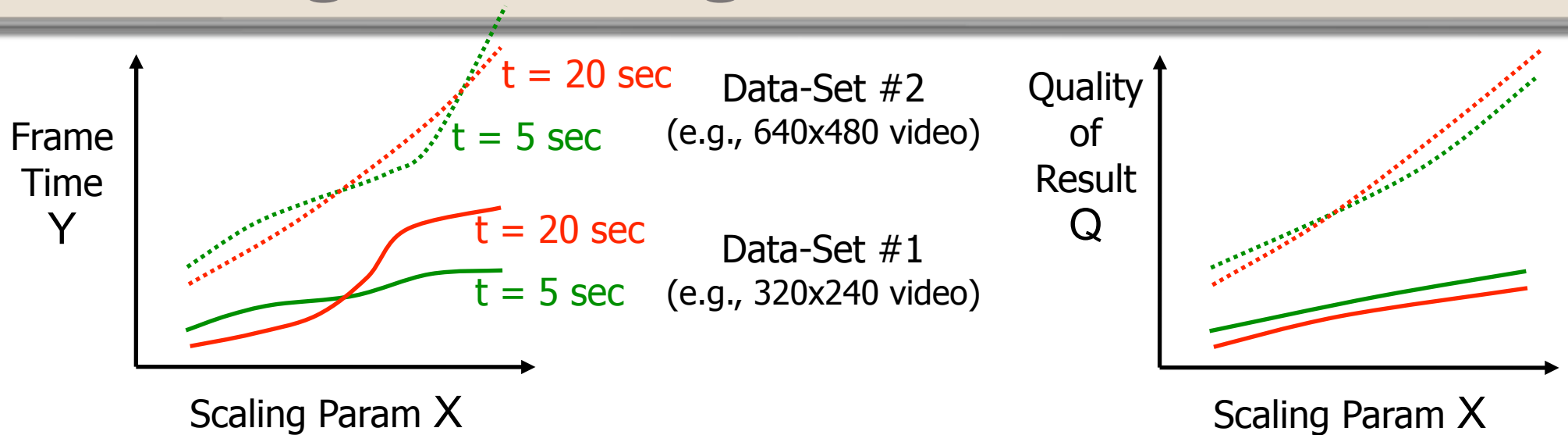
# Challenges in Tuning Feature Set



**X-Y**, **X-Q** Response Characteristics

- Highly time-varying, data-set dependent " *No fixed relationships*

- For practical purposes: UNKNOWN, entirely EMERGENT behavior

## So how do programmers currently tune?

- Impose *severe limitations*:
    - Video surveillance (slower motion, fixed background), at 320x240 resolution
    - Then: Fix X to *least harmful* value

# Challenges in Tuning Feature Set



**X-Y, X-Q** Response Characteristics

- Highly time-varying, data-set dependent " *No fixed relationships*

- For practical purposes: UNKNOWN, entirely EMERGENT behavior

## So how do programmers currently tune?

- Impose *severe limitations*:
  - Video surveillance (slower motion, fixed background), at 320x240 resolution
  - Then: Fix **X** to *least harmful* value

- Or, manually tune to *each game-play scenario*, tune for Xbox vs PS3

# Our Approach

# Our Approach

- Premise: Modeling emergent response behavior is hard!

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts


- Place only an *EASY* burden on the Programmer:
  - Identifies application-specific scale parameter X (WHAT to scale)

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

- Place only an *EASY* burden on the Programmer:
  - Identifies application-specific scale parameter $X$ (WHAT to scale)
  - Specifies desired frame-rate $Y$ to achieve (WHAT to achieve)

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

- Place only an *EASY* burden on the Programmer:
  - Identifies application-specific scale parameter $X$ (WHAT to scale)
  - Specifies desired frame-rate $Y$ to achieve (WHAT to achieve)

    (Implicit: Quality-of-Result $Q$ scales with $Y$)

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

- Place only an *EASY* burden on the Programmer:
  - Identifies application-specific scale parameter $X$ (WHAT to scale)
  - Specifies desired frame-rate $Y$ to achieve (WHAT to achieve)

    (Implicit: Quality-of-Result $Q$ scales with $Y$)

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

- Place only an *EASY* burden on the Programmer:
  - Identifies application-specific scale parameter X (WHAT to scale)
  - Specifies desired frame-rate Y to achieve (WHAT to achieve)

    (Implicit: Quality-of-Result Q scales with Y)

- HOW to achieve?

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

- Place only an *EASY* burden on the Programmer:
  - Identifies application-specific scale parameter $X$ (WHAT to scale)
  - Specifies desired frame-rate $Y$ to achieve (WHAT to achieve)

    (Implicit: Quality-of-Result $Q$ scales with $Y$)

- HOW to achieve?
  - *HARD*: Usually beyond scope of most programmers and even domain experts
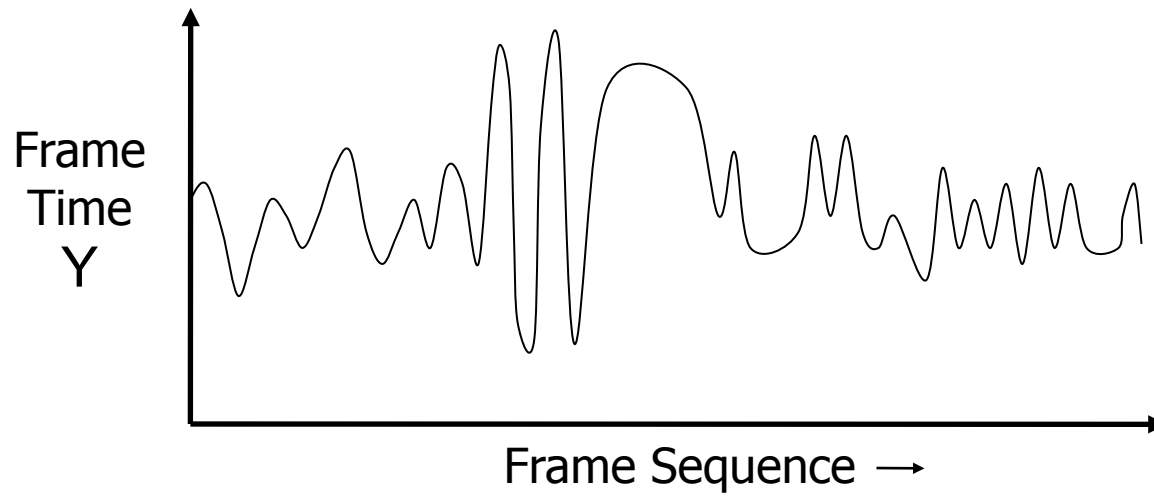
# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

- Place only an *EASY* burden on the Programmer:
  - Identifies application-specific scale parameter $X$ (WHAT to scale)
  - Specifies desired frame-rate $Y$ to achieve (WHAT to achieve)

    (Implicit: Quality-of-Result $Q$ scales with $Y$)

- HOW to achieve?
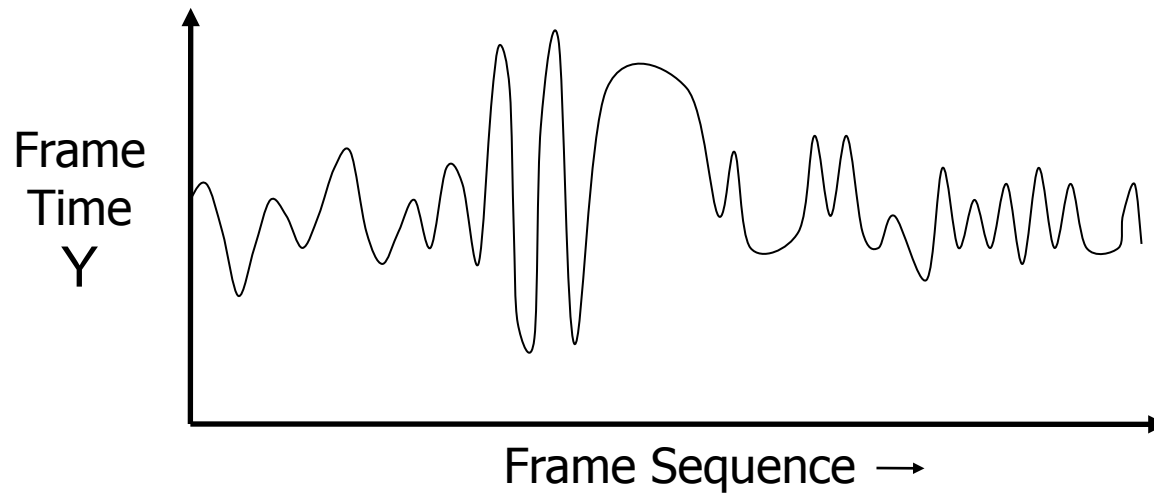  - *HARD*: Usually beyond scope of most programmers and even domain experts

  *What will it take for a **generic, broadly applicable Runtime Controller** to effectively control frame-QoS?*

# Our Approach

- Premise: Modeling emergent response behavior is hard!
  - Infeasible for programmers and even domain experts

- Place only an *EASY* burden on the Programmer:
  - Identifies application-specific scale parameter $X$ (WHAT to scale)
  - Specifies desired frame-rate $Y$ to achieve (WHAT to achieve)

    (Implicit: Quality-of-Result $Q$ scales with $Y$)

  Main Contribution
  Simplicity of Use,
  Generality of Application

- HOW to achieve?
  - *HARD*: Usually beyond scope of most programmers and even domain experts

  *What will it take for a **generic, broadly applicable Runtime Controller** to effectively control frame-QoS?*

# Control Problem

# Control Problem
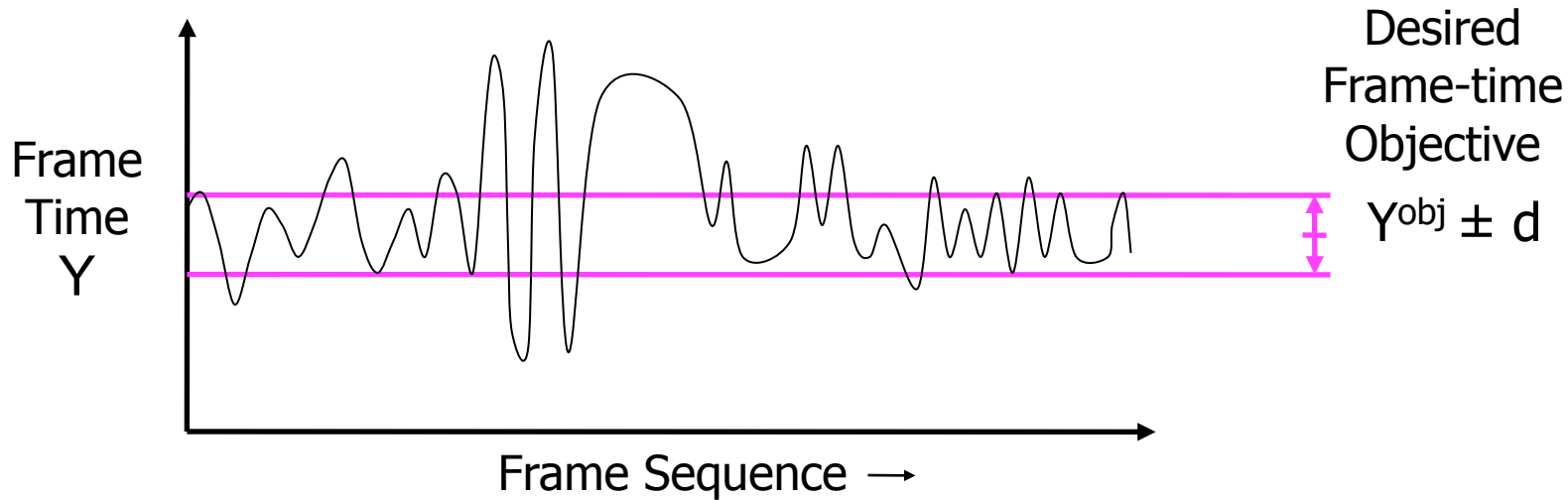


Frame Time Y vs. Frame Sequence →

- Soft Real Time nature
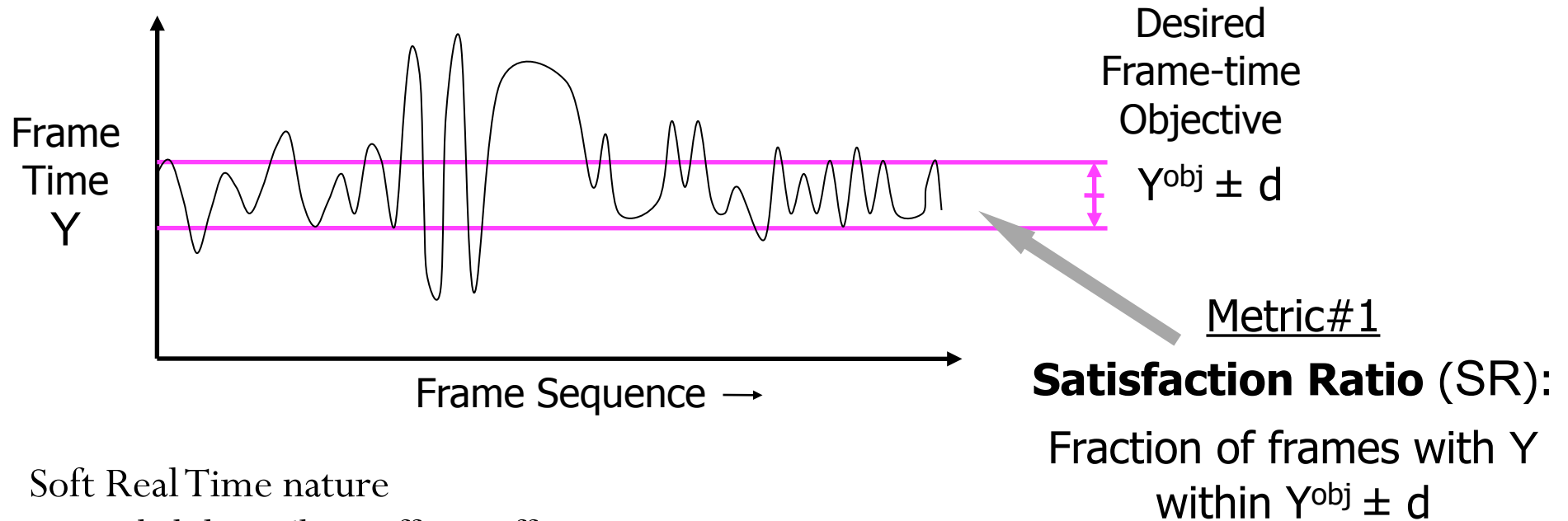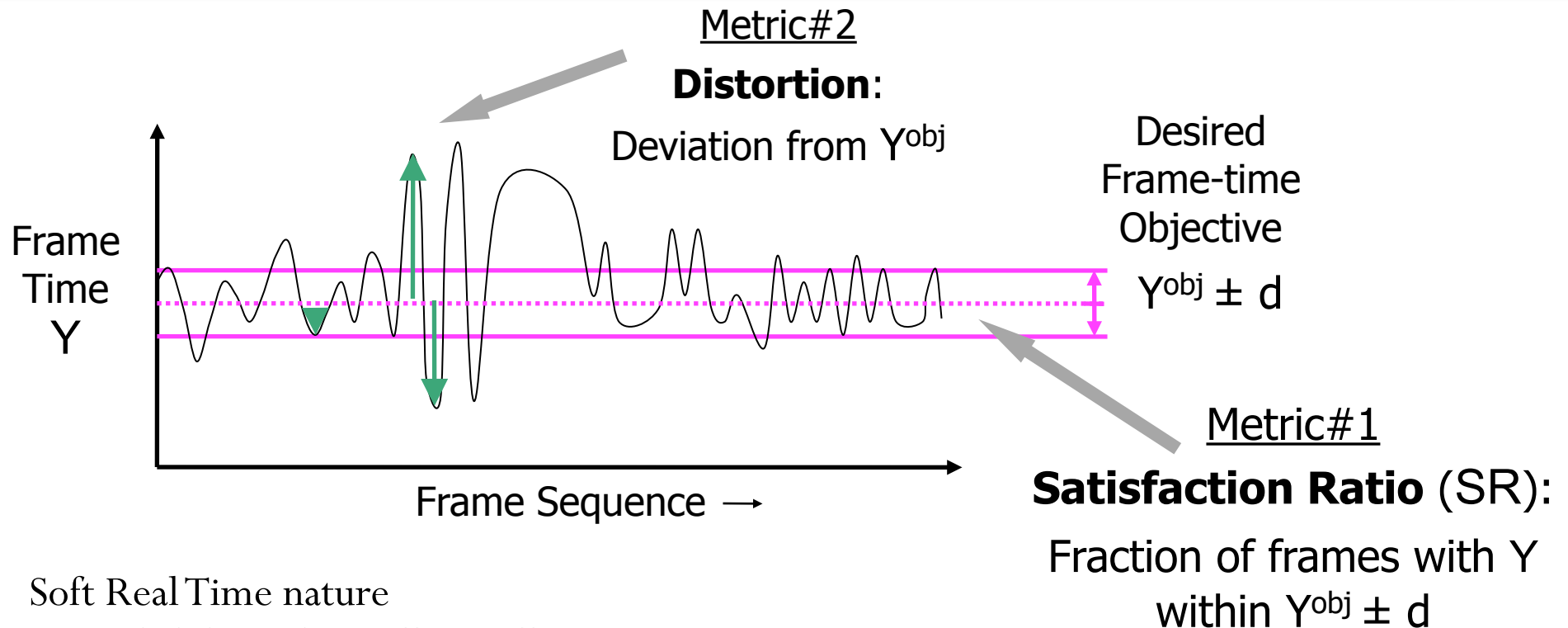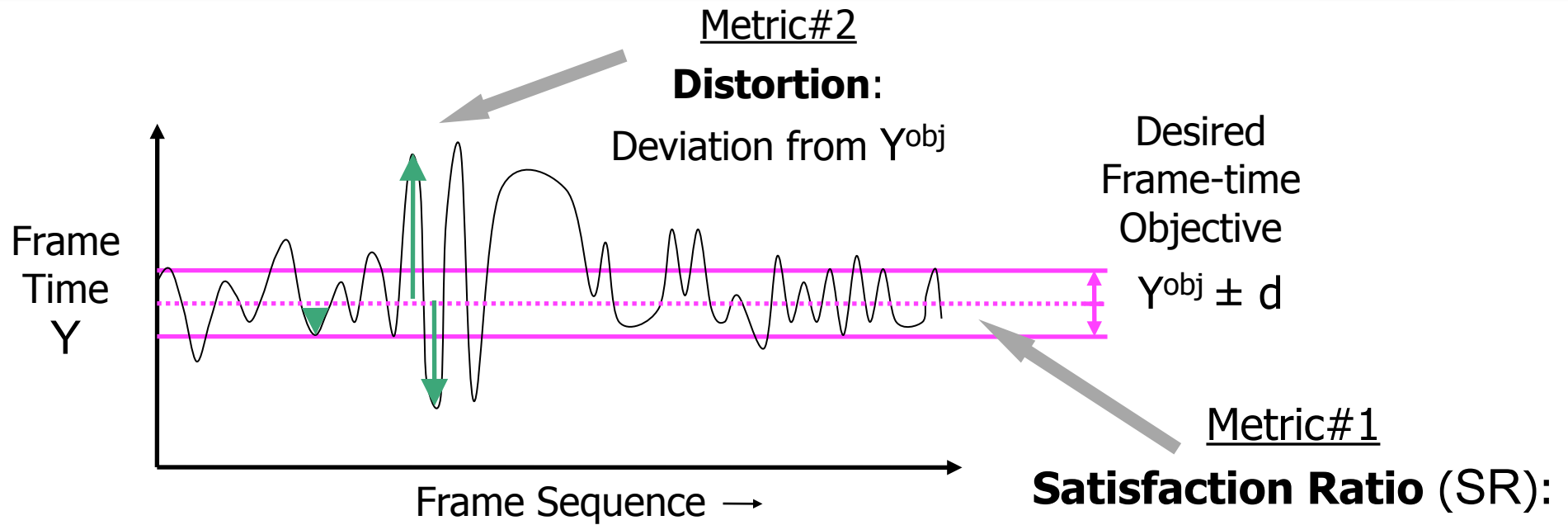
# Control Problem



- Soft Real Time nature
  - Probabilistic/best-effort suffices
    "Maintain 25-35 frames-per-second with high probability"

# Control Problem



- Soft Real Time nature
  - Probabilistic/best-effort suffices
    "Maintain 25-35 frames-per-second with high probability"

# Control Problem



- Soft Real Time nature
  - Probabilistic/best-effort suffices
    "Maintain 25-35 frames-per-second with high probability"

# Control Problem



- Soft Real Time nature
  - Probabilistic/best-effort suffices
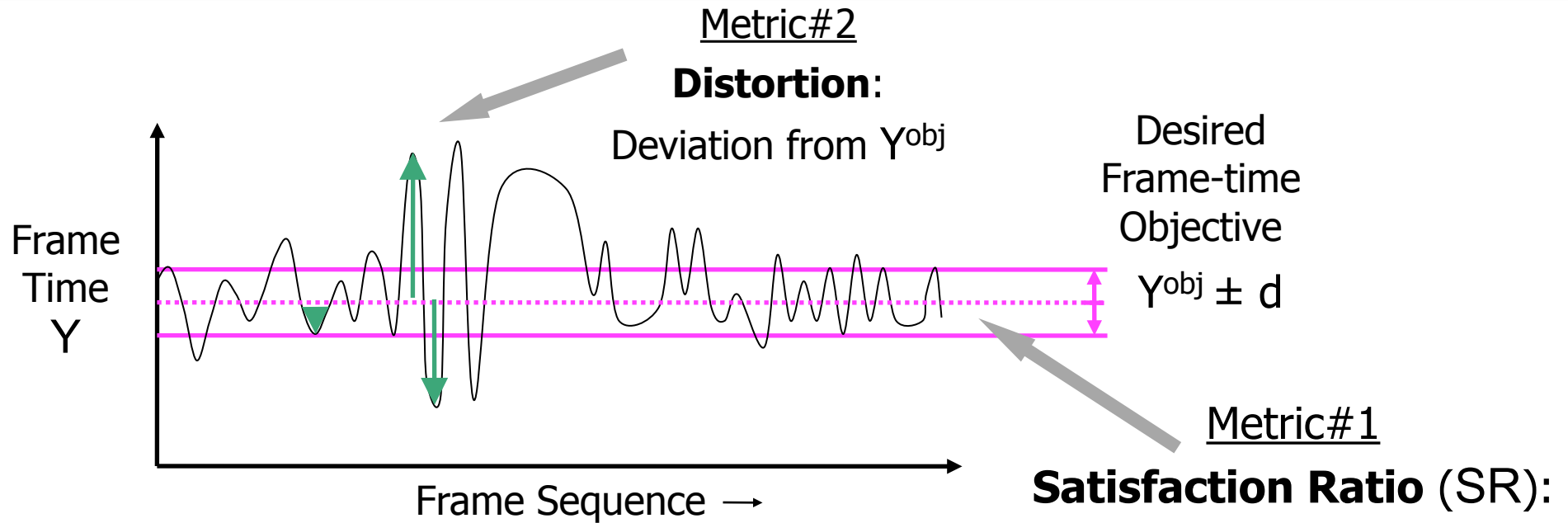    "Maintain 25-35 frames-per-second with high probability"

# Control Problem



Metric#2

**Distortion**:

Deviation from $Y^{obj}$

Desired
Frame-time
Objective

$Y^{obj} \pm d$

Frame
Time
Y

Frame Sequence →

Metric#1

**Satisfaction Ratio** (SR):

Fraction of frames with Y
within $Y^{obj} \pm d$

- Soft Real Time nature
  - Probabilistic/best-effort suffices
    "Maintain 25-35 frames-per-second with high probability"
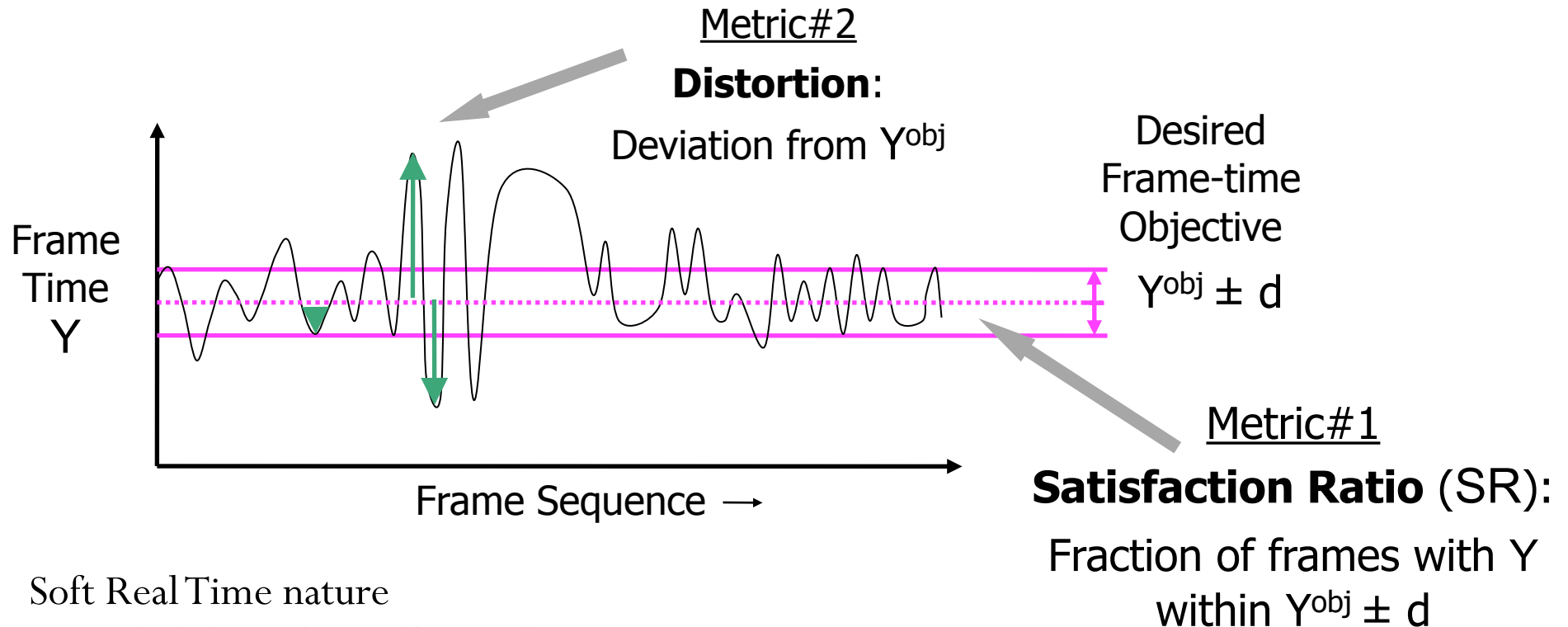
- What's a good **SR** on a data-set?

# Control Problem



- Soft Real Time nature
  - Probabilistic/best-effort suffices
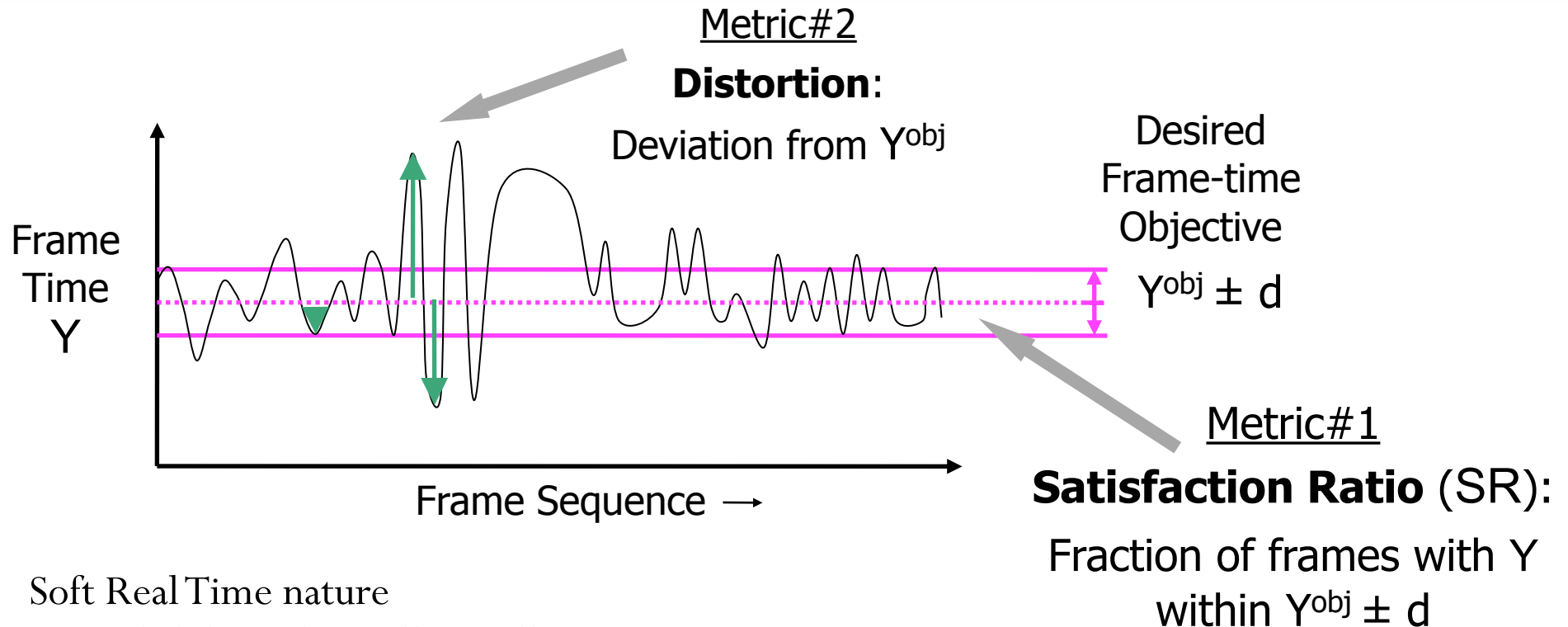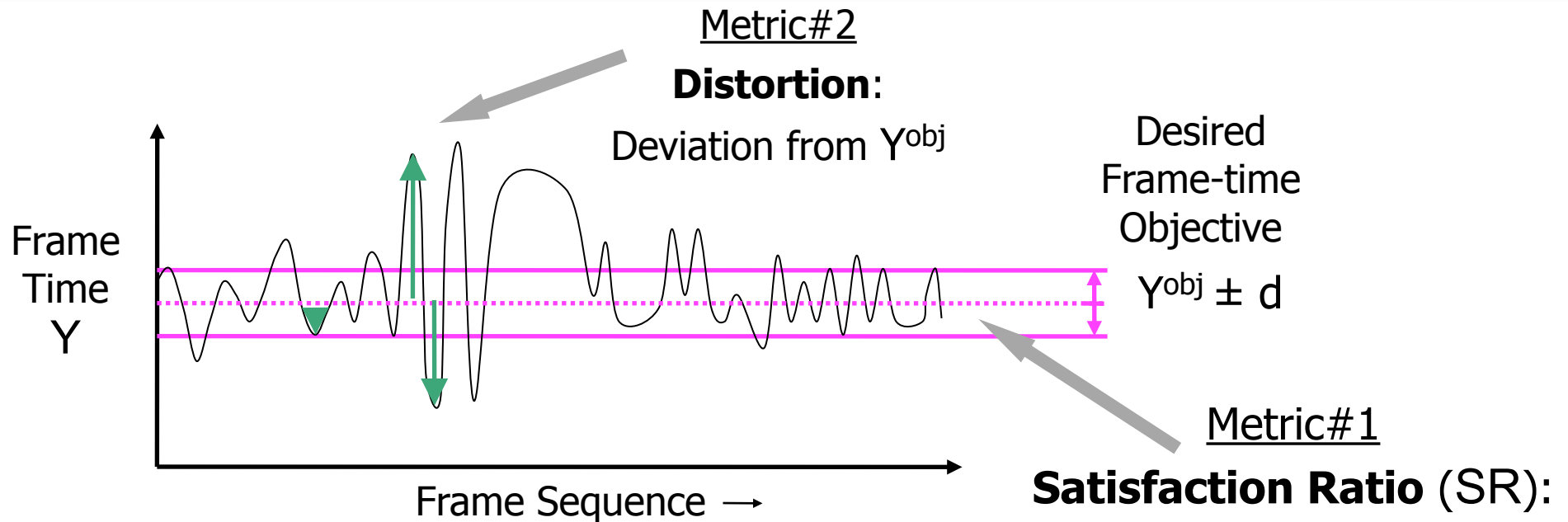    "Maintain 25-35 frames-per-second with high probability"

- What's a good **SR** on a data-set?
  - *30%, 50%, 90%, 99%?*

# Control Problem



**Metric#2**

**Distortion**:

Deviation from $Y^{obj}$

Desired
Frame-time
Objective

$Y^{obj} \pm d$

Frame
Time
Y

Frame Sequence →

**Metric#1**

**Satisfaction Ratio** (SR):

Fraction of frames with Y
within $Y^{obj} \pm d$

- Soft Real Time nature
  - Probabilistic/best-effort suffices
    "Maintain 25-35 frames-per-second with high probability"

- What's a good SR on a data-set?
  - *30%, 50%, 90%, 99%?*
  - A specified objective $Y^{obj} \pm d$ may be *impossible* to achieve

# Control Problem

Metric#2

**Distortion**:

Deviation from $Y^{obj}$

Desired
Frame-time
Objective

$Y^{obj} \pm d$

Frame
Time
$Y$

Frame Sequence →

Metric#1

**Satisfaction Ratio** (SR):

Fraction of frames with $Y$
within $Y^{obj} \pm d$

- Soft Real Time nature
  - Probabilistic/best-effort suffices
    "Maintain 25-35 frames-per-second with high probability"

- What's a good **SR** on a data-set?
  - *30%, 50%, 90%, 99%?*
  - A specified objective $Y^{obj} \pm d$ may be ***impossible*** to achieve
  - In general, application characteristics are unknown, \ best feasible **SR** is unknown

# Control Problem



**Metric#2**
**Distortion**:
Deviation from $Y^{obj}$

Desired Frame-time Objective
$Y^{obj} \pm d$

Frame Time Y

Frame Sequence →

**Metric#1**
**Satisfaction Ratio** (SR):
Fraction of frames with Y within $Y^{obj} \pm d$

- Soft Real Time nature
  - Probabilistic/best-effort suffices
    "Maintain 25-35 frames-per-second with high probability"

- What's a good **SR** on a data-set?
  - *30%, 50%, 90%, 99%?*
  - A specified objective $Y^{obj} \pm d$ may be *impossible* to achieve
  - In general, application characteristics are unknown, \ best feasible **SR** is unknown
  " Only option: compare against **SR** when **X** has *best fixed setting*

# Runtime Controller: Attempt#1

- Goal: Find **best fixed** $X$ for each data-set
  - Highly Sub-Optimal!

# Runtime Controller: Attempt#1

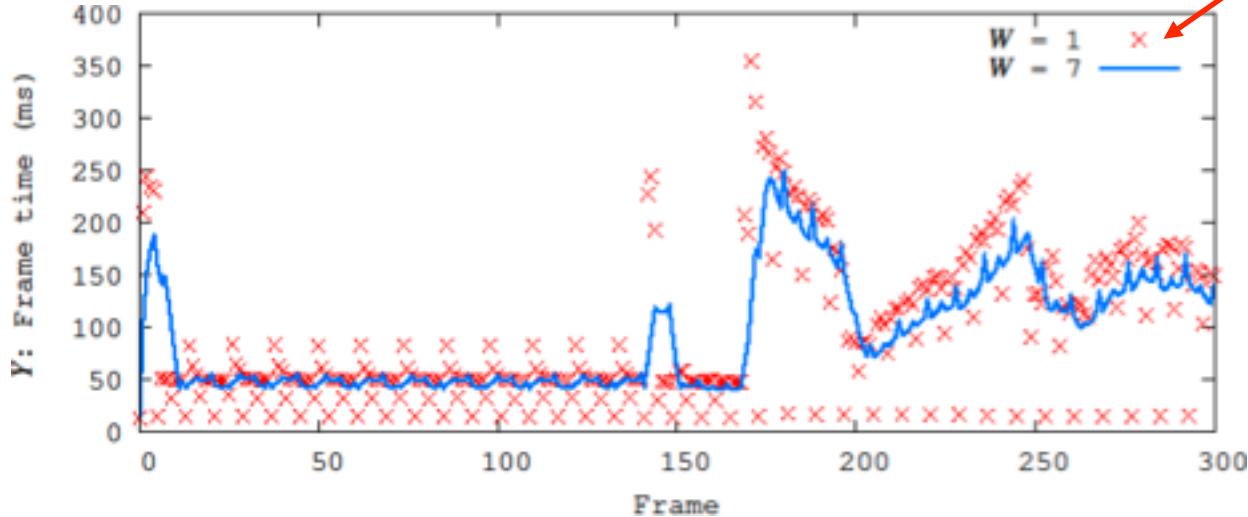- Goal: Find *best fixed* X for each data-set
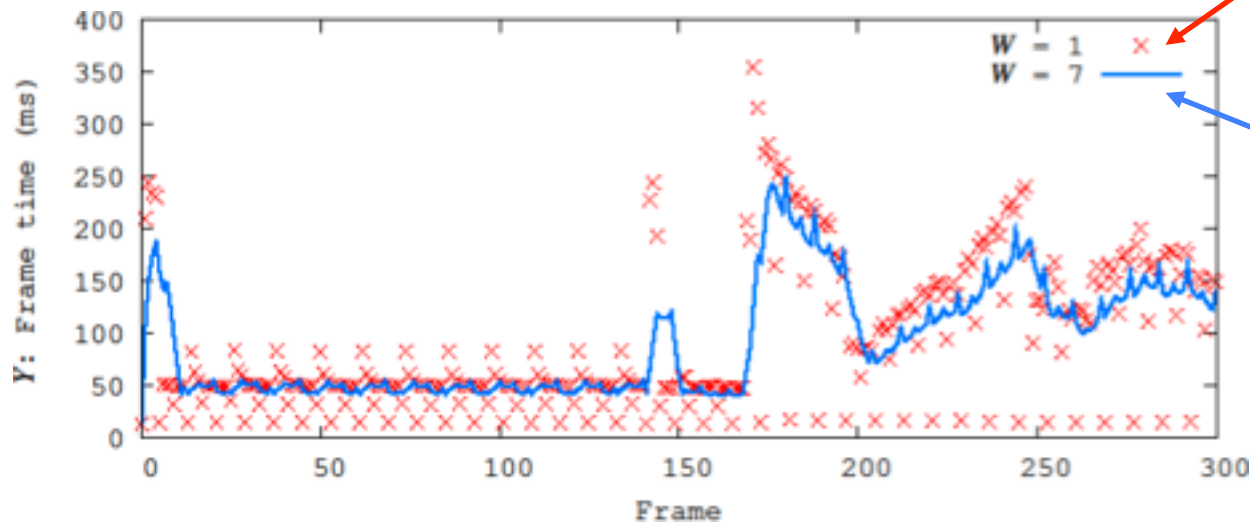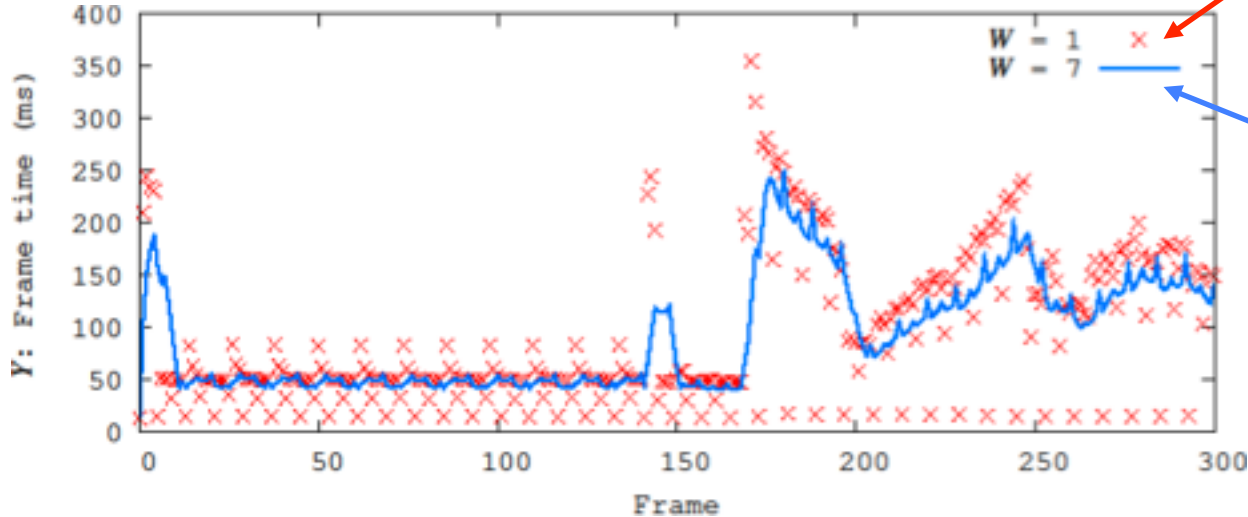  - Highly Sub-Optimal!

### MPEG2 Encoder Frame Sequence (Fixed X)

- Goal: Find **best fixed** X for each data-set
  - Highly Sub-Optimal!

MPEG2 Encoder Frame Sequence (Fixed X)

Instantaneous frame-times have significant transients *(Not perceived by user!)*

# Runtime Controller: Attempt#1

- Goal: Find **best fixed** X for each data-set
  - Highly Sub-Optimal!

## MPEG2 Encoder Frame Sequence (Fixed X)



Instantaneous frame-times have significant transients *(Not perceived by user!)*

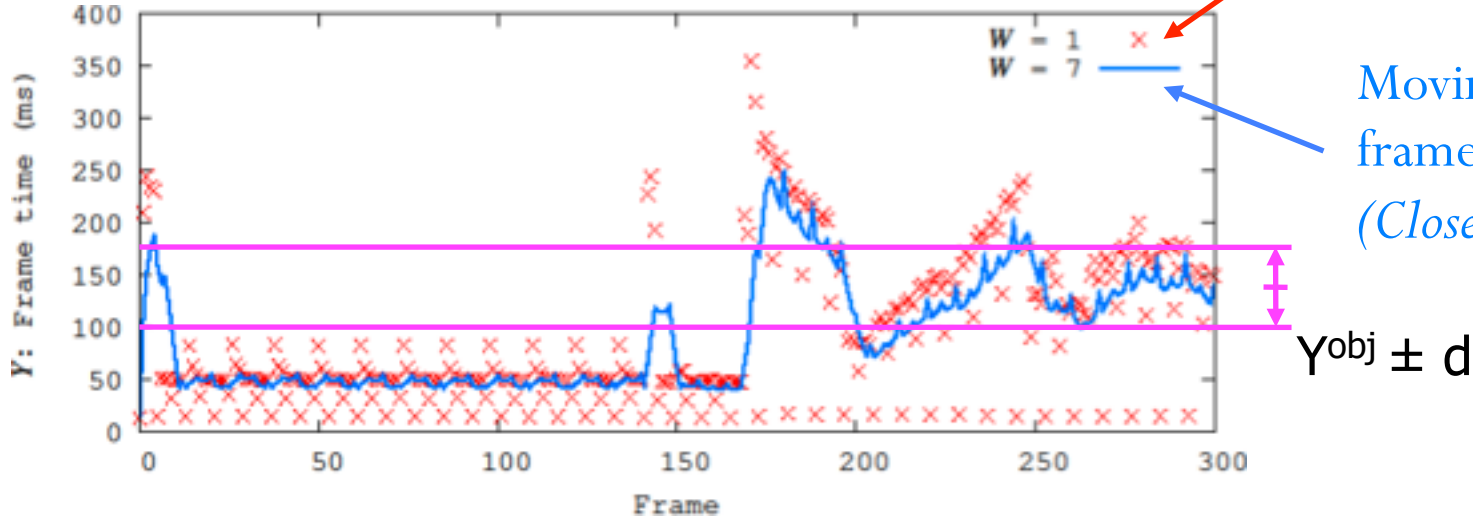Moving average of previous 7 frames is mostly smooth *(Closer to user perception)*

# Runtime Controller: Attempt#1

- Goal: Find *best fixed* X for each data-set
  - Highly Sub-Optimal!

MPEG2 Encoder Frame Sequence (Fixed X)



Instantaneous frame-times have significant transients *(Not perceived by user!)*

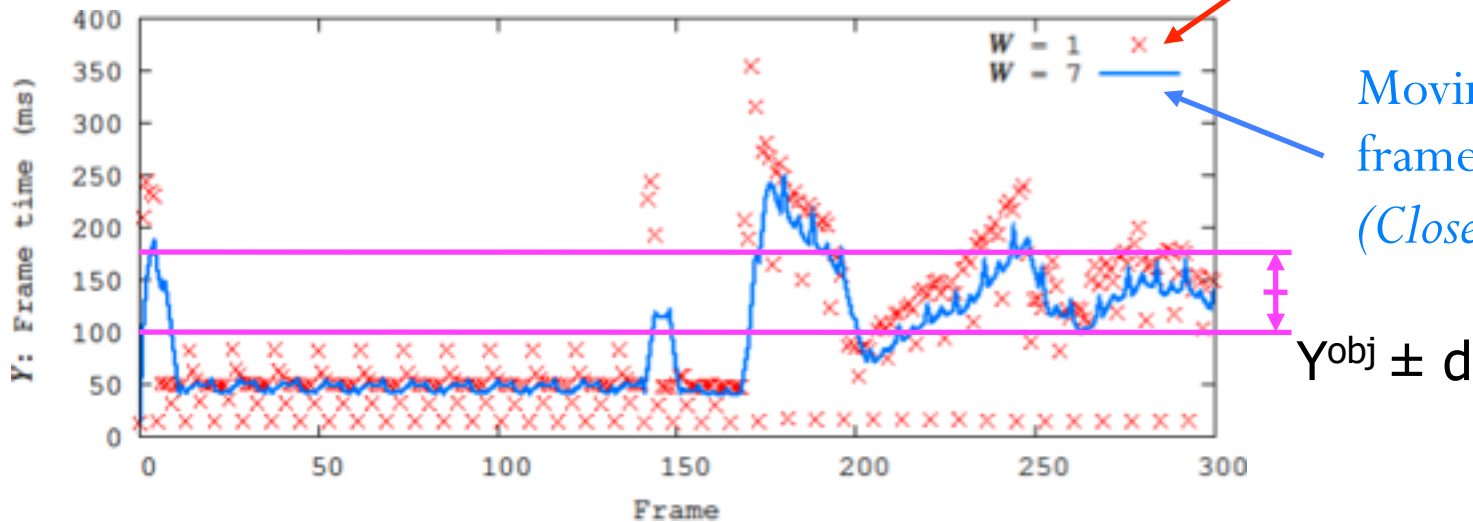Moving average of previous 7 frames is mostly smooth *(Closer to user perception)*

Lesson#1

-Control **perceived** frame-time (much easier)

-Programmer specifies: Window of Perception W

7

# Runtime Controller: Attempt#1

- Goal: Find *best fixed* X for each data-set
  - Highly Sub-Optimal!

MPEG2 Encoder Frame Sequence (Fixed X)



Instantaneous frame-times have significant transients *(Not perceived by user!)*

Moving average of previous 7 frames is mostly smooth *(Closer to user perception)*
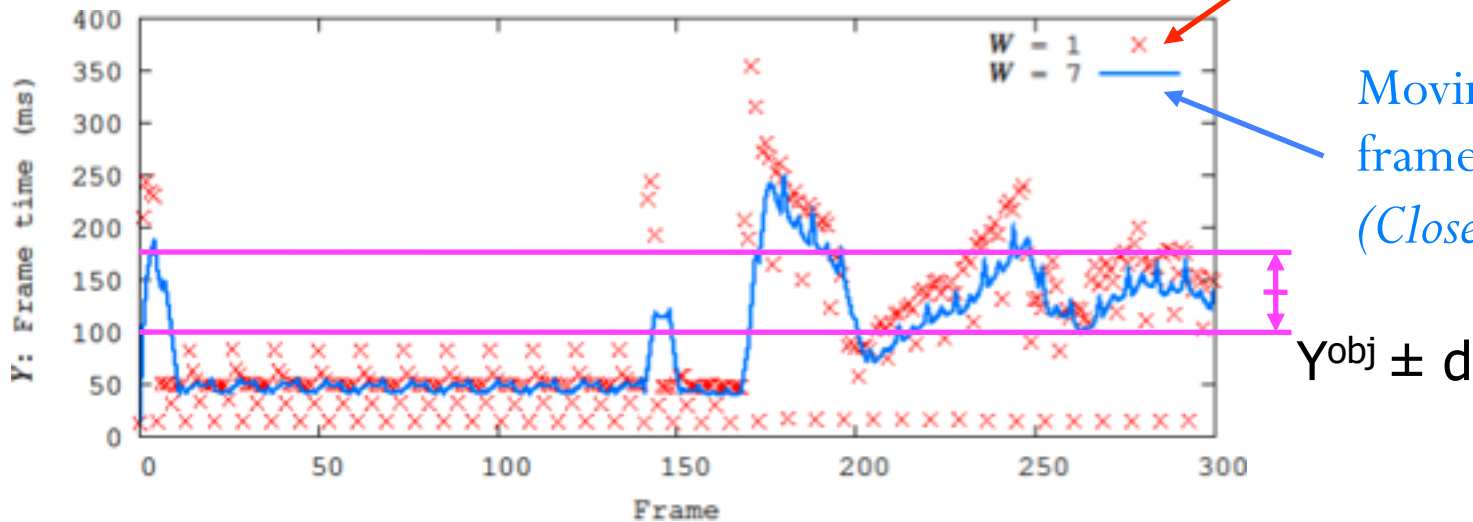
$Y^{obj} \pm d$

Lesson#1

-Control **perceived** frame-time (much easier)

-Programmer specifies: Window of Perception W

# Runtime Controller: Attempt#1

- Goal: Find *best fixed* X for each data-set
  - Highly Sub-Optimal!

**MPEG2 Encoder Frame Sequence (Fixed X)**



Instantaneous frame-times have significant transients *(Not perceived by user!)*

Moving average of previous 7 frames is mostly smooth *(Closer to user perception)*

$Y^{obj} \pm d$

Execution characteristics can vary significantly over regions
(**Regional Dependence**)

Lesson#1
-Control **perceived** frame-time (much easier)
-Programmer specifies: Window of Perception W

7

# Runtime Controller: Attempt#1

- Goal: Find **best fixed** X for each data-set
  - Highly Sub-Optimal!

Instantaneous frame-times have significant transients *(Not perceived by user!)*

MPEG2 Encoder Frame Sequence (Fixed X)

Moving average of previous 7 frames is mostly smooth *(Closer to user perception)*



$Y^{obj} \pm d$

Execution characteristics can vary significantly over regions

(**Regional Dependence**)

Lesson#2
Train X to each region!

Lesson#1
-Control **perceived** frame-time (much easier)
-Programmer specifies: Window of Perception W

# Runtime Controller: Attempt#2

# Runtime Controller: Attempt#2

- Goal: Train X to each region

# Runtime Controller: Attempt#2

- Goal: Train **X** to each region

- Reinforcement learning of **X-Y** relationship
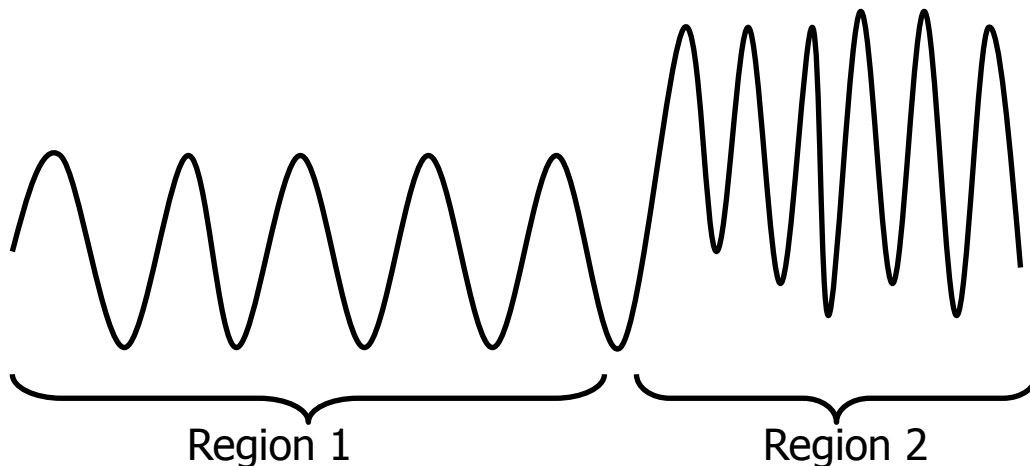  - Too slow!

# Runtime Controller: Attempt#2

- Goal: Train **X** to each region

- Reinforcement learning of **X-Y** relationship
  - Too slow!
    **X-Y** relationships in MPEG2 Encoder and Torque Game usually change before learning can be used to subsequently control Y

- Find good **X** for region before **X-Y** relationship changes significantly

# Runtime Controller: Attempt#2

- Goal: Train **X** to each region

- Reinforcement learning of **X-Y** relationship
  - Too slow!
    **X-Y** relationships in MPEG2 Encoder and Torque Game usually change before learning can be used to subsequently control Y

- Find good **X** for region before **X-Y** relationship changes significantly
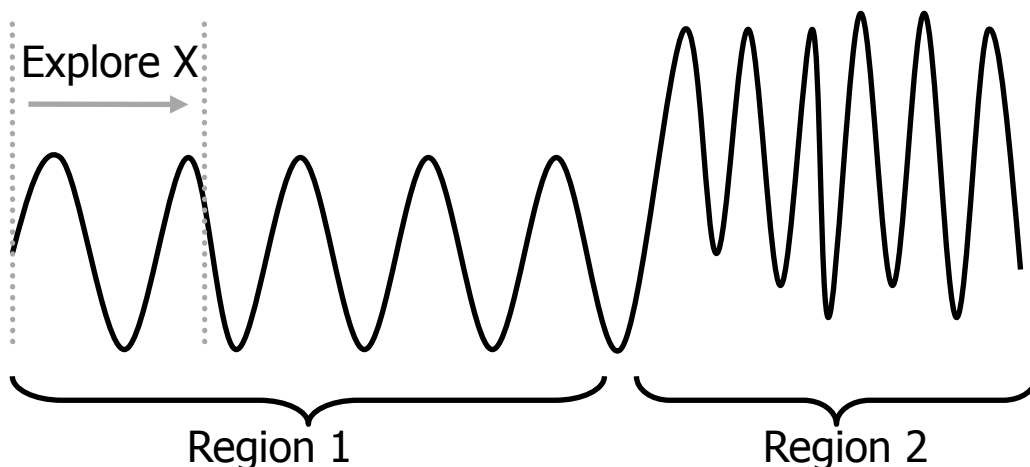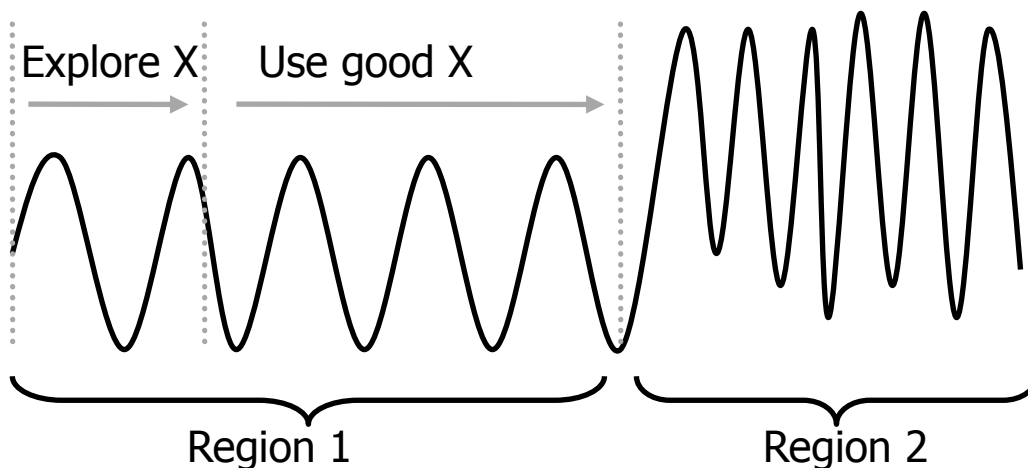- Then exploit good **X** for rest of (hopefully long) region " achieves high **SR**

- Goal: Train **X** to each region

- Reinforcement learning of **X-Y** relationship
  - Too slow!
    **X-Y** relationships in MPEG2 Encoder and Torque Game usually change before learning can be used to subsequently control Y

- Find good **X** for region before **X-Y** relationship changes significantly
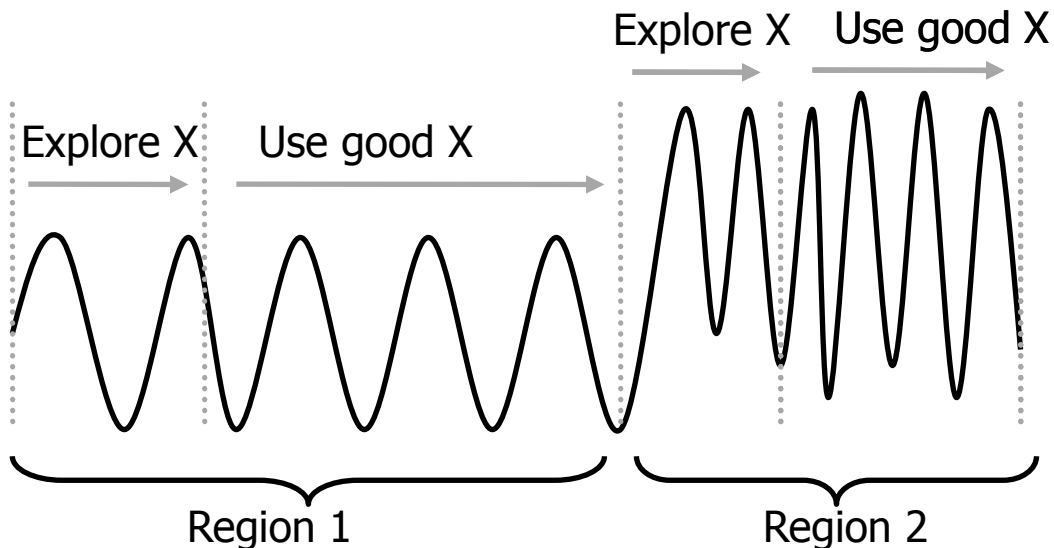- Then exploit good **X** for rest of (hopefully long) region " achieves high **SR**



Region 1          Region 2

# Runtime Controller: Attempt#2

- Goal: Train **X** to each region

- Reinforcement learning of **X-Y** relationship
  - Too slow!
    **X-Y** relationships in MPEG2 Encoder and Torque Game usually change before learning can be used to subsequently control Y

- Find good **X** for region before **X-Y** relationship changes significantly
- Then exploit good **X** for rest of (hopefully long) region " achieves high **SR**

Explore X

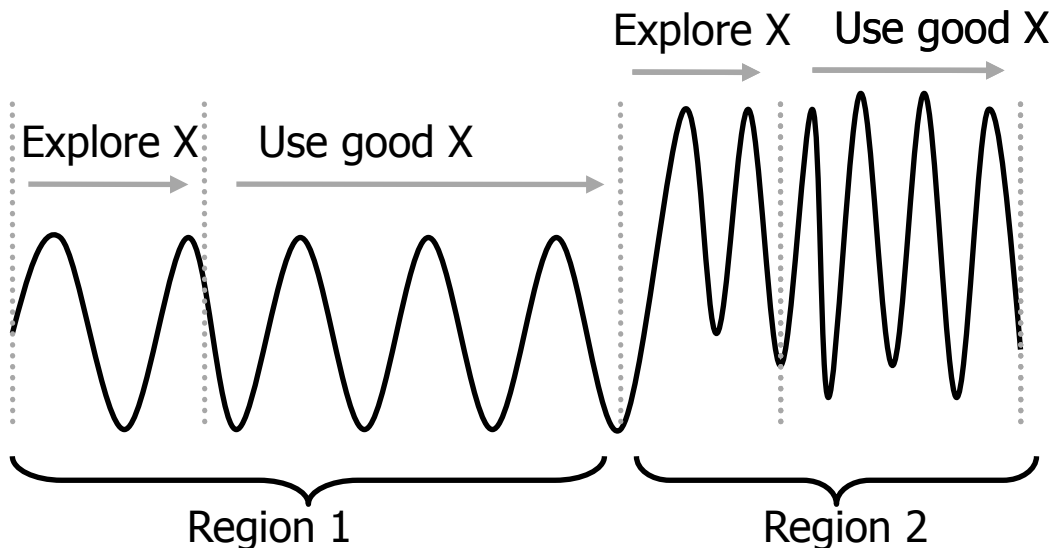Region 1

Region 2

# Runtime Controller: Attempt#2

- Goal: Train **X** to each region

- Reinforcement learning of **X-Y** relationship
  - Too slow!
    **X-Y** relationships in MPEG2 Encoder and Torque Game usually change before learning can be used to subsequently control Y

- Find good **X** for region before **X-Y** relationship changes significantly
- Then exploit good **X** for rest of (hopefully long) region " achieves high **SR**

Explore X    Use good X

Region 1                 Region 2

# Runtime Controller: Attempt#2

- Goal: Train **X** to each region

- Reinforcement learning of **X-Y** relationship
  - Too slow!
    **X-Y** relationships in MPEG2 Encoder and Torque Game usually change before learning can be used to subsequently control Y

- Find good **X** for region before **X-Y** relationship changes significantly
- Then exploit good **X** for rest of (hopefully long) region " achieves high **SR**

Explore X    Use good X

Explore X    Use good X

Region 1                    Region 2

# Runtime Controller: Attempt#2

- Goal: Train **X** to each region

- Reinforcement learning of **X-Y** relationship
  - Too slow!
    **X-Y** relationships in MPEG2 Encoder and Torque Game usually change before learning can be used to subsequently control Y

- Find good **X** for region before **X-Y** relationship changes significantly
- Then exploit good **X** for rest of (hopefully long) region " achieves high **SR**

Explore X    Use good X

Explore X    Use good X

Region 1                Region 2

Observation

Application "world state" does not change dramatically faster than user perception time
Lesson#3

Stable Region Length >> W, with high probability

# Runtime Controller: Attempt#3

# Runtime Controller: Attempt#3

- Use Feedback Controller: adjust **X** based on observed error $\Delta Y$

# Runtime Controller: Attempt#3

- Use Feedback Controller: adjust $X$ based on observed error $\Delta Y$
  - Need to assume *monotonicity* in $X$-$Y$ relationship

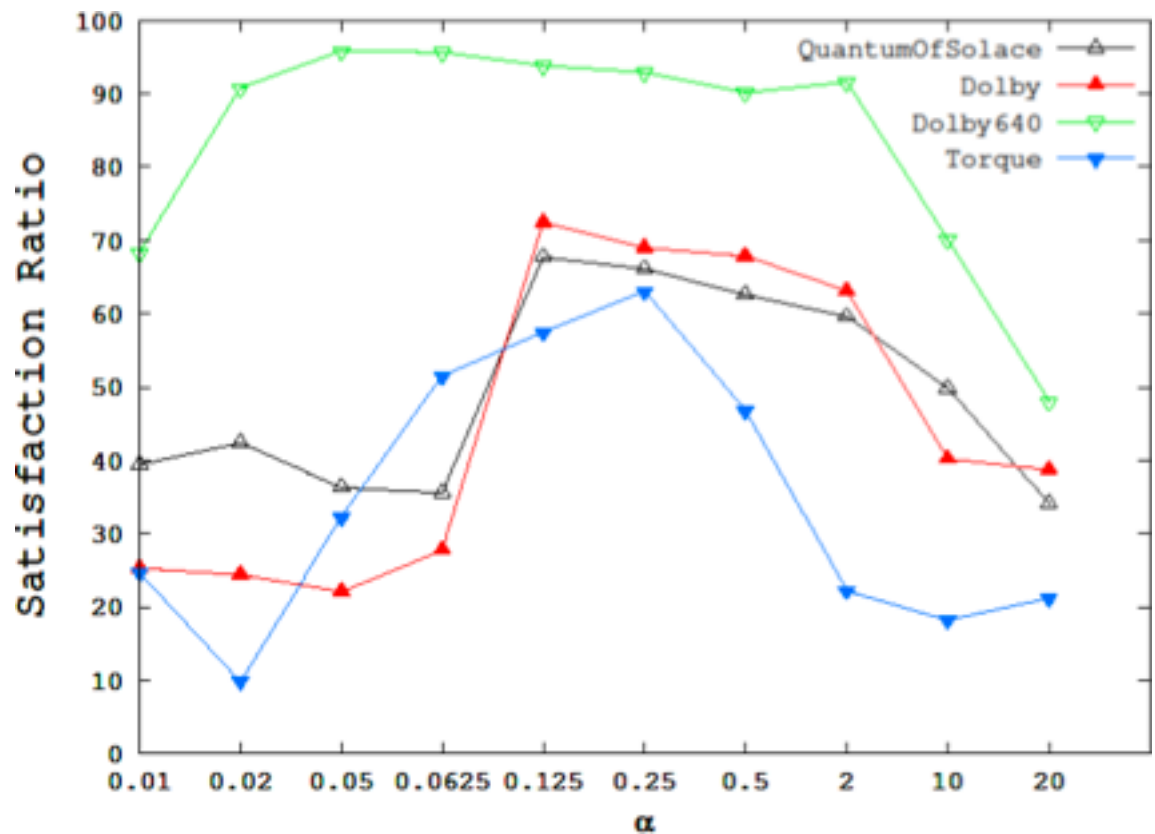# Runtime Controller: Attempt#3

- Use Feedback Controller: adjust $X$ based on observed error $\Delta Y$
  - Need to assume *monotonicity* in $X$-$Y$ relationship
  - But $X$-$Y$ relationship can remain time-varying

- Simplest: P controller

# Runtime Controller: Attempt#3

- Use Feedback Controller: adjust **X** based on observed error $\Delta Y$
  - Need to assume *monotonicity* in **X-Y** relationship
  - But **X-Y** relationship can remain time-varying

- Simplest: P controller

$$\Delta X \leftarrow (1/\alpha) \; \Delta Y$$

# Runtime Controller: Attempt#3

- Use Feedback Controller: adjust **X** based on observed error $\Delta Y$

  - Need to assume ***monotonicity*** in **X-Y** relationship

  - But **X-Y** relationship can remain time-varying

- Simplest: P controller

  $\Delta X \leftarrow (1/\alpha)\ \Delta Y$

  ( Controller Gain: **1/α**)

# Runtime Controller: Attempt#3

- Use Feedback Controller: adjust **X** based on observed error $\Delta Y$
  - Need to assume **_monotonicity_** in **X-Y** relationship
  - But **X-Y** relationship can remain time-varying

- Simplest: P controller

  $\Delta X \leftarrow (1/\alpha) \, \Delta Y$

  ( Controller Gain: **1/α**)



9

# Runtime Controller: Attempt#3

- Use Feedback Controller: adjust **X** based on observed error **ΔY**
  - Need to assume ***monotonicity*** in **X-Y** relationship
  - But **X-Y** relationship can remain time-varying

- Simplest: P controller

  $$\Delta X \leftarrow (1/\alpha)\ \Delta Y$$
  ( Controller Gain: **1/α** )

**Lesson#3**
Feedback Control works
poorly, unless $\alpha$ is in a
narrow range determined
by data-set

# Adaptive Feedback Controller

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  #1. Monotonic $X$-$Y$

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  #1. Monotonic $X-Y$
  - Makes rapid feedback-control feasible, dramatically simplifies learning

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  - #1. Monotonic $X$-$Y$
    - Makes rapid feedback-control feasible, dramatically simplifies learning

  - #2. Only perceived frame-rate matters, $W$ specified

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller

  #1. Monotonic $X$-$Y$
  - Makes rapid feedback-control feasible, dramatically simplifies learning

  #2. Only perceived frame-rate matters, $W$ specified
  - Minimizes influence of frequent transients

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  #1. Monotonic $X$-$Y$
  - Makes rapid feedback-control feasible, dramatically simplifies learning

  #2. Only perceived frame-rate matters, $W$ specified
  - Minimizes influence of frequent transients
  - When is "failure" too long? → Perceptible to user

  #3. Stable Region length $>> W$

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  - #1. Monotonic $X\text{-}Y$
    - Makes rapid feedback-control feasible, dramatically simplifies learning

  - #2. Only perceived frame-rate matters, $W$ specified
    - Minimizes influence of frequent transients
    - When is "failure" too long? $\rightarrow$ Perceptible to user

  - #3. Stable Region length $>> W$
    - Determines quantitative metrics and tests to adapt policy

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  - #1. Monotonic $X$-$Y$
    - Makes rapid feedback-control feasible, dramatically simplifies learning

  - #2. Only perceived frame-rate matters, $W$ specified
    - Minimizes influence of frequent transients
    - When is "failure" too long? $\rightarrow$ Perceptible to user

  - #3. Stable Region length $>> W$
    - Determines quantitative metrics and tests to adapt policy

- Simplicity in Feedback Law:

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller

  #1. Monotonic $X$-$Y$
  - Makes rapid feedback-control feasible, dramatically simplifies learning

  #2. Only perceived frame-rate matters, $W$ specified
  - Minimizes influence of frequent transients
  - When is "failure" too long? → Perceptible to user

  #3. Stable Region length $>> W$
  - Determines quantitative metrics and tests to adapt policy

- Simplicity in Feedback Law:
  $$\Delta X \leftarrow (1/\alpha)\, \Delta Y$$

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  - #1. Monotonic $X$-$Y$
    - Makes rapid feedback-control feasible, dramatically simplifies learning

  - #2. Only perceived frame-rate matters, $W$ specified
    - Minimizes influence of frequent transients
    - When is "failure" too long? → Perceptible to user

  - #3. Stable Region length $>> W$
    - Determines quantitative metrics and tests to adapt policy

- Simplicity in Feedback Law:
  $$\Delta X \leftarrow (1/\alpha) \, \Delta Y$$
- Sophistication in Adaptive Policy:

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  - #1. Monotonic $X\text{-}Y$
    - Makes rapid feedback-control feasible, dramatically simplifies learning

  - #2. Only perceived frame-rate matters, $W$ specified
    - Minimizes influence of frequent transients
    - When is "failure" too long? $\rightarrow$ Perceptible to user

  - #3. Stable Region length $>> W$
    - Determines quantitative metrics and tests to adapt policy

- Simplicity in Feedback Law:
  $$\Delta X \leftarrow (1/\alpha)\ \Delta Y$$
- Sophistication in Adaptive Policy:
  How/When to adjust $\alpha$?

# Adaptive Feedback Controller

- Programmer specifies: $X$, $Y^{obj} \pm d$, $W$
- Objective: maintain a high $SR$

- **Domain Assumptions** help us design generic controller
  #1. Monotonic $X$-$Y$
    - Makes rapid feedback-control feasible, dramatically simplifies learning

  #2. Only perceived frame-rate matters, $W$ specified
    - Minimizes influence of frequent transients
    - When is "failure" too long? → Perceptible to user

  #3. Stable Region length $>> W$
    - Determines quantitative metrics and tests to adapt policy

- Simplicity in Feedback Law:
  $\Delta X \leftarrow (1/\alpha) \, \Delta Y$
- Sophistication in Adaptive Policy:
  How/When to adjust $\alpha$?

Sufficient to deliver significant improvements in SR

# Adaptive Policy Design

# Adaptive Policy Design

- Criteria for **Significant Policy failure**
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*

# Adaptive Policy Design

- Criteria for **Significant Policy failure**
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

# Adaptive Policy Design

- Criteria for ***Significant Policy failure***
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

- Strategy

# Adaptive Policy Design

- Criteria for **_Significant Policy failure_**
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

- Strategy
  - P controller's simplicity '' clearly identifiable **Failure Modes**

# Adaptive Policy Design

- Criteria for *Significant Policy failure*
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

- Strategy
  - P controller's simplicity " clearly identifiable **Failure Modes**
  - $\alpha$ needs to be corrected within narrow data-set dependent range (for high $SR$)

# Adaptive Policy Design

- Criteria for ***Significant Policy failure***
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

- Strategy
  - P controller's simplicity '' clearly identifiable **Failure Modes**
  - $\alpha$ needs to be corrected within narrow data-set dependent range (for high $SR$)

  **Global Failure** (e.g., 320x240 vs 640x480 data-sets):
  *orders-of-magnitude correction in $\alpha$ in single step*
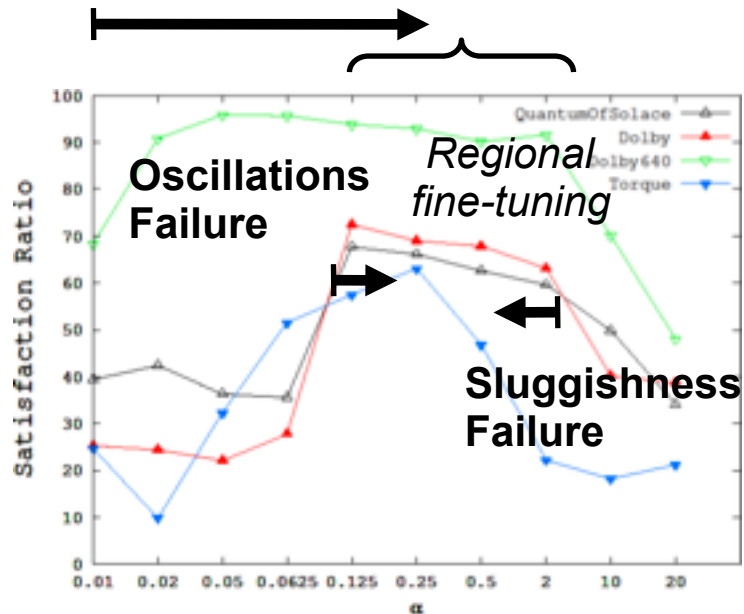
# Adaptive Policy Design

- Criteria for *Significant Policy failure*
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

- Strategy
  - P controller's simplicity " clearly identifiable **Failure Modes**
  - $\alpha$ needs to be corrected within narrow data-set dependent range (for high $SR$)

  **Global Failure** (e.g., 320x240 vs 640x480 data-sets):
  *orders-of-magnitude correction in α in single step*

# Adaptive Policy Design

- Criteria for *Significant Policy failure*
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

- Strategy
  - P controller's simplicity '' clearly identifiable **Failure Modes**
  - $\alpha$ needs to be corrected within narrow data-set dependent range (for high $SR$)

**Global Failure** (e.g., 320x240 vs 640x480 data-sets):
*orders-of-magnitude correction in α in single step*

# Adaptive Policy Design

- Criteria for *Significant Policy failure*
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

- Strategy
  - P controller's simplicity '' clearly identifiable **Failure Modes**
  - $\alpha$ needs to be corrected within narrow data-set dependent range (for high $SR$)

  **Global Failure** (e.g., 320x240 vs 640x480 data-sets):
  *orders-of-magnitude correction in α in single step*

# Adaptive Policy Design

- Criteria for *Significant Policy failure*
  - $\Delta Y$ compared to $Y^{obj}$, and compared to $2*d$ *(error of excessive magnitude)*
  - $Y$ outside $Y^{obj} \pm d$ for more than $W$ frames *(error of excessive duration)*

- Strategy
  - P controller's simplicity '' clearly identifiable **Failure Modes**
  - $\alpha$ needs to be corrected within narrow data-set dependent range (for high $SR$)

**Global Failure** (e.g., 320x240 vs 640x480 data-sets):
*orders-of-magnitude correction in α in single step*

# Illustration of Oscillation Failure Mode

# Illustration of Oscillation Failure Mode

- Failure Metrics
  - H, L
  - $\eta \tilde{} \, d * \eta + H * W / L, \ (0 < d < 1)$

# Illustration of Oscillation Failure Mode

- Failure Metrics
  - H, L
  - $\eta \tilde{} \ d * \eta + H * W / L, \ (0 < d < 1)$

- Failure Detection
  - $\eta > t$
    with $t = 1/(1-d) * 2d * 1.0$

- Policy Adaptation
  - $\alpha^{new} \tilde{} \ \alpha * \eta / t$

# Illustration of Oscillation Failure Mode

- Failure Metrics
  - H, L
  - $\eta \sim d * \eta + H * W / L,\ (0 < d < 1)$

- Failure Detection
  - $\eta > t$
    with t = 1/(1-d) * 2d * 1.0

- Policy Adaptation
  - $\alpha^{new} \sim \alpha * \eta / t$

- *Accumulative metrics* for failure
  - Weighted to gradually forget old learning, at rate d
  - Robustness against transients
  - Responsive to persistent changed behavior
  - Constant state " very low runtime overhead of controller

# Execution Trace: MPEG2 Encoder



$X \rightarrow Search\,Window\,Size : \{0 \rightarrow 30, 1 \rightarrow 20, 2 \rightarrow 15, 3 \rightarrow 10, 4 \rightarrow 5, 5 \rightarrow 2, 6 \rightarrow 1, 7 \rightarrow 0\}$

**Integral X**: Easier for programmer to just **sample** *Search Window Size* over sufficient range

# Benchmark Result: MPEG2 Encoder

# Benchmark Result: MPEG2 Encoder



Dolby

- Video sequence: dolbycity320x240
- Adaptive Controller better than *envelope of best fixed* Xs
  - Due to Regional tuning of X vs only global tuning

# Benchmark Result: MPEG2 Encoder

# Benchmark Result: MPEG2 Encoder



- Video sequence: dolbycity 640x480
- Adaptive better than every Fixed X case overall
  - Even though for a given $Y^{obj}$, a particular Fixed X might match or exceed Adaptive

# Benchmark Result: Torque Game Engine



- Quality of Result: *Gameplay Intelligence*

# Benchmark Result: Torque Game Engine



Torque

- Quality of Result: *Gameplay Intelligence*
- For $Y^{obj} = 0.04$ secs
  - Adaptive: 24ms of AI/frame
  - Best fixed X: 14ms of AI/frame

# Related Work

# Related Work

- Real-time techniques already address QoS and Soft Real-time. BUT:
  - Application needs to be implemented as Task-Graphs
  - With execution-time properties specified for nodes [Mejia-Alvarez et al. RTSS 2000]
  - And, Utility functions for QoS provided [Block et al. ECRTS 2008]

- Application-specific techniques
  - Rate-distortion control in H.264 encoder
  - QoS control in MPEG2 *decoder* (easier than *encoder* due to video-control-sequence information) [Roitzsch et al. RTSS'06] [Huang et al. MULTIMEDIA'07] [Wust et al. ECRTS'04]

# Related Work

- Real-time techniques already address QoS and Soft Real-time. BUT:
    - Application needs to be implemented as Task-Graphs
    - With execution-time properties specified for nodes [Mejia-Alvarez et al. RTSS 2000]
    - And, Utility functions for QoS provided [Block et al. ECRTS 2008]

- Application-specific techniques
    - Rate-distortion control in H.264 encoder
    - QoS control in MPEG2 *decoder* (easier than *encoder* due to video-control-sequence information) [Roitzsch et al. RTSS'06] [Huang et al. MULTIMEDIA'07] [Wust et al. ECRTS'04]

- "Ready-made" Adaptive techniques
    - Illustrated by: Web-cache QoS optimization [Lu et al. IWQoS 2002]
    
      Adaptive Pole-placement controller-design, based on periodically re-estimating a $2^{nd}$ order parametric LTI model

# Related Work

- Real-time techniques already address QoS and Soft Real-time. BUT:
  - Application needs to be implemented as Task-Graphs
  - With execution-time properties specified for nodes [Mejia-Alvarez et al. RTSS 2000]
  - And, Utility functions for QoS provided [Block et al. ECRTS 2008]

- Application-specific techniques
  - Rate-distortion control in H.264 encoder
  - QoS control in MPEG2 *decoder* (easier than *encoder* due to video-control-sequence information) [Roitzsch et al. RTSS'06] [Huang et al. MULTIMEDIA'07] [Wust et al. ECRTS'04]

- "Ready-made" Adaptive techniques
  - Illustrated by: Web-cache QoS optimization [Lu et al. IWQoS 2002]

    Adaptive Pole-placement controller-design, based on periodically re-estimating a $2^{nd}$ order parametric LTI model

- Our work is an *Adaptive Gain-Scheduling* controller
  - Adjust control-law directly to overcome observed failures

# Related Work

- Real-time techniques already address QoS and Soft Real-time. BUT:
  - Application needs to be implemented as Task-Graphs
  - With execution-time properties specified for nodes [Mejia-Alvarez et al. RTSS 2000]
  - And, Utility functions for QoS provided [Block et al. ECRTS 2008]

- Application-specific techniques
  - Rate-distortion control in H.264 encoder
  - QoS control in MPEG2 *decoder* (easier than *encoder* due to video-control-sequence information) [Roitzsch et al. RTSS'06] [Huang et al. MULTIMEDIA'07] [Wust et al. ECRTS'04]

- "Ready-made" Adaptive techniques
  - Illustrated by: Web-cache QoS optimization [Lu et al. IWQoS 2002]

    Adaptive Pole-placement controller-design, based on periodically re-estimating a 2nd order parametric LTI model

- Our work is an *Adaptive Gain-Scheduling* controller
  - Adjust control-law directly to overcome observed failures
- Controller is *domain-specific* rather than *application-specific*
  - Based on broad domain-assumptions, rather than linear dynamical models of plant
  - *Semi-physical System-Identification*

# Related Work

- Real-time techniques already address QoS and Soft Real-time. BUT:
  - Application needs to be implemented as Task-Graphs
  - With execution-time properties specified for nodes [Mejia-Alvarez et al. RTSS 2000]
  - And, Utility functions for QoS provided [Block et al. ECRTS 2008]

- Application-specific techniques
  - Rate-distortion control in H.264 encoder
  - QoS control in MPEG2 *decoder* (easier than *encoder* due to video-control-sequence information) [Roitzsch et al. RTSS'06] [Huang et al. MULTIMEDIA'07] [Wust et al. ECRTS'04]

- "Ready-made" Adaptive techniques
  - Illustrated by: Web-cache QoS optimization [Lu et al. IWQoS 2002]
    Adaptive Pole-placement controller-design, based on periodically re-estimating a $2^{nd}$ order parametric LTI model

- Our work is an *Adaptive Gain-Scheduling* controller
  - Adjust control-law directly to overcome observed failures
- Controller is *domain-specific* rather than *application-specific*
  - Based on broad domain-assumptions, rather than linear dynamical models of plant
  - *Semi-physical System-Identification*
- Very effective, but Best-Effort:
  - No guarantees on Safety, Robustness, Reachability

# Conclusion

# Conclusion

- Our controller greatly simplifies frame-QoS control
  - Programmers avoid having to model emergent behaviors
  - Incorporate controller as a library
    - exceedingly light-weight: < 0.05% overhead

# Conclusion

- Our controller greatly simplifies frame-QoS control
  - Programmers avoid having to model emergent behaviors
  - Incorporate controller as a library
    - exceedingly light-weight: < 0.05% overhead

- Generality
  - Any application that satisfies Domain Assumptions
    - Monotonic response, perceived frame-rate, Stable response periods >> W frames
  - Graceful degradation when not

- No general-purpose alternatives available that work well for interactive applications whose behavior is *UNKNOWN*, *EMERGENT* and *NON-ANALYZABLE*

# Conclusion

- Our controller greatly simplifies frame-QoS control
  - Programmers avoid having to model emergent behaviors
  - Incorporate controller as a library
    - exceedingly light-weight: < 0.05% overhead

- Generality
  - Any application that satisfies Domain Assumptions
    - Monotonic response, perceived frame-rate, Stable response periods >> W frames
  - Graceful degradation when not

- No general-purpose alternatives available that work well for interactive applications whose behavior is *UNKNOWN*, *EMERGENT* and *NON-ANALYZABLE*

- Future Work
  - Multiple X, Multiple prioritized Y, Explicit Q
  - Domain Observations "Adaptive Control + Least-Squares Function Estimation
  - But, much more compute intensive controller

# Thank you!

- Questions?

# *Torque* Game Engine: Measured Behavior



*objective*: 25 to 42 fps

# *Torque* Game Engine: Measured Behavior

# *Torque* Game Engine: Measured Behavior

# Distortion in Torque

# Standard Deviation in Torque

**Games, Multimedia, Interactive Viz**

# Don't Real-Time Methods Solve This Already?

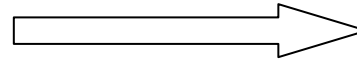**Games, Multimedia, Interactive Viz** ⟶
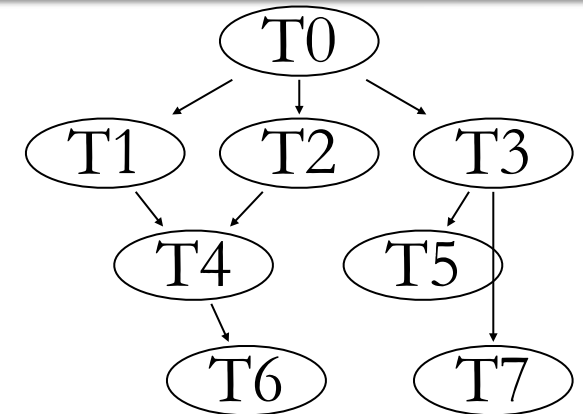
**Implement as a Real-Time App**

**Games, Multimedia, Interactive Viz**

**Implement as a Real-Time App**

# Don't Real-Time Methods Solve This Already?

**Games, Multimedia, Interactive Viz**

$\longrightarrow$



**Implement as a Real-Time App**

**Real-Time Task-Graph**
- Application decomposed into Tasks and Precedence Constraints

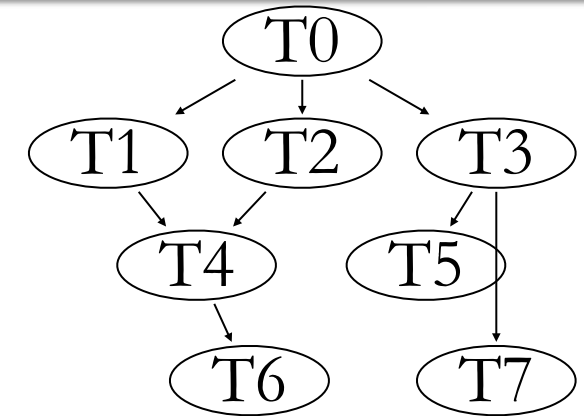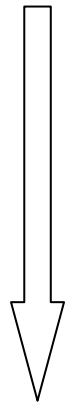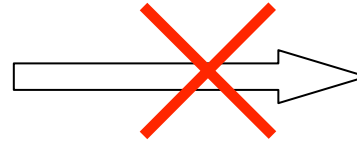- Responsiveness guaranteed by Real-time semantics (hard or probabilistic)

# Games, Multimedia, Interactive Viz

**Implement as a Real-Time App**

**Implement with High-Productivity, Large Scale Programming flows**

T0

T1  T2  T3

T4  T5

T6  T7

**Real-Time Task-Graph**
- Application decomposed into Tasks and Precedence Constraints

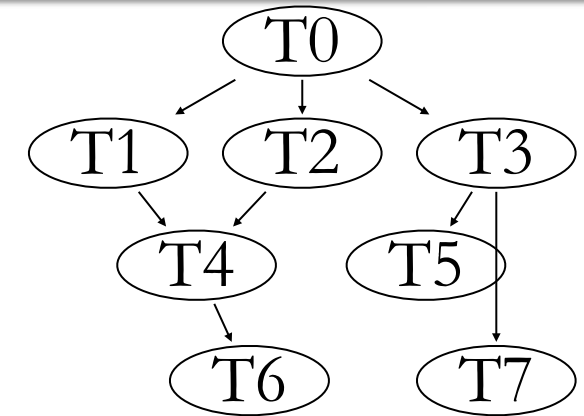- Responsiveness guaranteed by Real-time semantics (hard or probabilistic)

# Don't Real-Time Methods Solve This Already?

**Games, Multimedia, Interactive Viz**



**Implement as a Real-Time App**
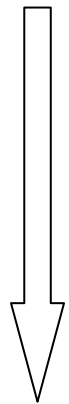
**Implement with High-Productivity, Large Scale Programming flows**

**Real-Time Task-Graph**

- Application decomposed into Tasks and Precedence Constraints

- Responsiveness guaranteed by Real-time semantics (hard or probabilistic)

# Don't Real-Time Methods Solve This Already?

**Games, Multimedia, Interactive Viz**



**Implement as a Real-Time App**

**Implement with High-Productivity, Large Scale Programming flows**

**Real-Time Task-Graph**

- Application decomposed into Tasks and Precedence Constraints

- Responsiveness guaranteed by Real-time semantics (hard or probabilistic)
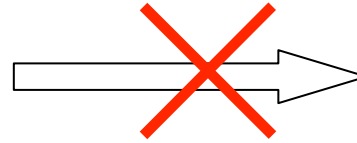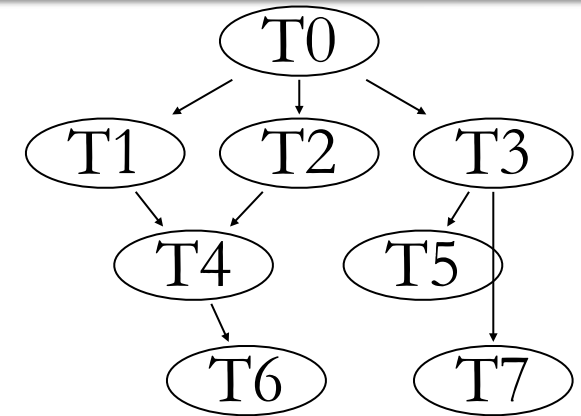
**C, C++, Java: Monolithic App**

- 100Ks to Millions of LoC

- No analyzable structure for responsiveness and scaling

- Responsiveness and Quality entirely **emergent** attributes (currently tuning this is an art)

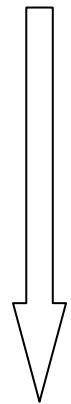# Don't Real-Time Methods Solve This Already?

## Games, Multimedia, Interactive Viz



**Implement as a Real-Time App**

**Implement with High-Productivity, Large Scale Programming flows**

**C, C++, Java: Monolithic App**

- 100Ks to Millions of LoC

- No analyzable structure for responsiveness and scaling

- Responsiveness and Quality entirely **emergent** attributes (currently tuning this is an art)

**Real-Time Task-Graph**

- Application decomposed into Tasks and Precedence Constraints

- Responsiveness guaranteed by Real-time semantics (hard or probabilistic)

**Need a new bag of tricks to Scale Semantics in Monolithic Applications**

23

# Runtime Controller