

Lecture 18:

Differentiable programming in Slang + Presentation Talk Tips

**Visual Computing Systems
Stanford CS348K, Spring 2025**

Two reasons to discuss Slang

- 1. Interesting programming system of potentially high practical utility**
- 2. Great example to discuss policy and mechanism in systems design**

Policy vs. mechanism in programming systems

- 1. Mechanism(s): tools for performing simple, primitive operations**
- 2. Policy: how the mechanisms are used to perform a higher-level task**

Let's consider some examples from class so far:

Halide

- Mechanisms for defining allocation of buffers and iteration order of computation loops
- Halide tries hard to not be opinionated about how developer should use those mechanisms (but some policy w.r.t. always implementing “sliding window optimizations” when they are possible.

Frankencamera:

- Mechanisms for controlling/configuring the sensor and other camera devices and for specifying the time when the operations should happen
- One major policy design in the system: “best effort” adherence to user's commands

Consider these systems with some popular web frameworks you might be familiar with (Django, Flask, Express.js), etc...

Why is it important to be intentional about separating policy from mechanism?

- **Generality of system: well-designed primitives could be used and composed in many ways**
- **Different application areas/tasks might demand different policies, and the choice of the right policy might be critical to success, with only the developer of the application able to know what the best policy is**
 - **Error handling policy: Applications might differ in whether they want a best-effort system or not (examples: Frankencamera, TCP vs. IP, etc.)**
 - **Scheduling policy: different image processing pipelines might demand very different scheduling strategies (fuse everything, vs burn memory but save recompute)**
 - **Policy of what algorithms to use: some programs might be amenable to forward model diff, others reverse mode**
- **Systems that stick to providing powerful mechanisms (and stay clear of dictating policy), and allow policy decisions to be implemented by layers above them, tend to be long-lived systems**

Slang

- **Programming language and compiler mechanisms to enable differentiable programming in a wide range of use cases...**
 - **But is carefully designed to avoid policy decisions as much as possible**
- **Key mechanism: ability to perform auto-differentiation of a “basic block” of code (but this is not unique to Slang, there are many differentiable programming systems)**
- **But Slang is unique in that it: anticipates a diverse set of application programs**
 - **Code that is most efficiently implements using forward mode autodiff or reverse-mode autodiff**
 - **Code that mixes differentiable and non-differentiable logic**
 - **Code containing functions that in isolation are not differentiable but programmer has additional knowledge about how they want the derivative computed (overrides)**

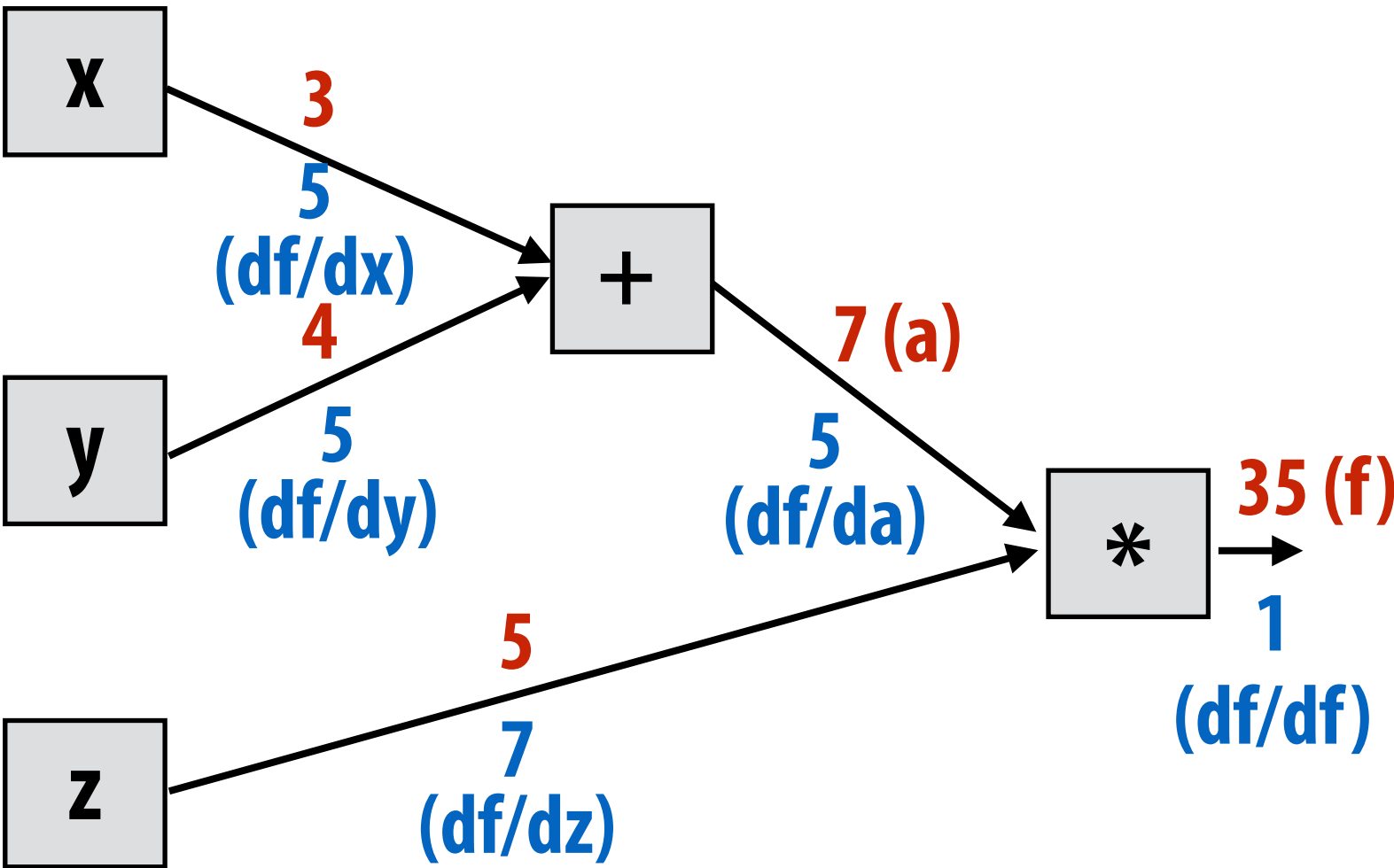
A bit of a reminder: chain rule

$f(x, y, z) = (x + y)z = az$ **Where:** $a = x + y$

$\frac{df}{da} = z$ $\frac{da}{dx} = 1$ $\frac{da}{dy} = 1$

So, by the derivative chain rule:

$\frac{df}{dx} = \frac{df}{da} \frac{da}{dx} = z$

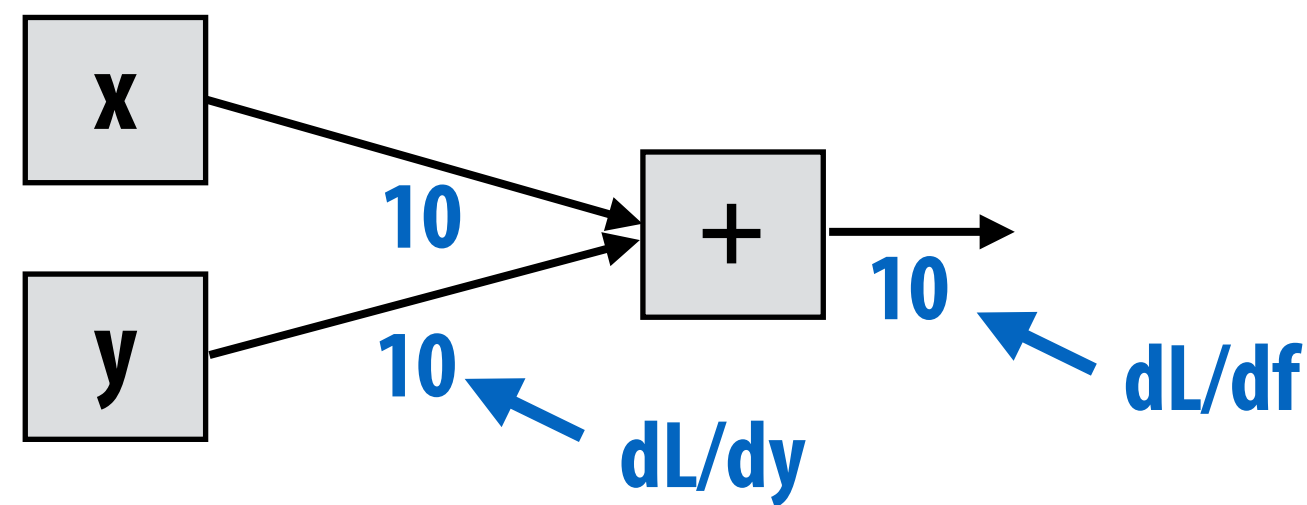


Red = output of node
Blue = $df/dnode$

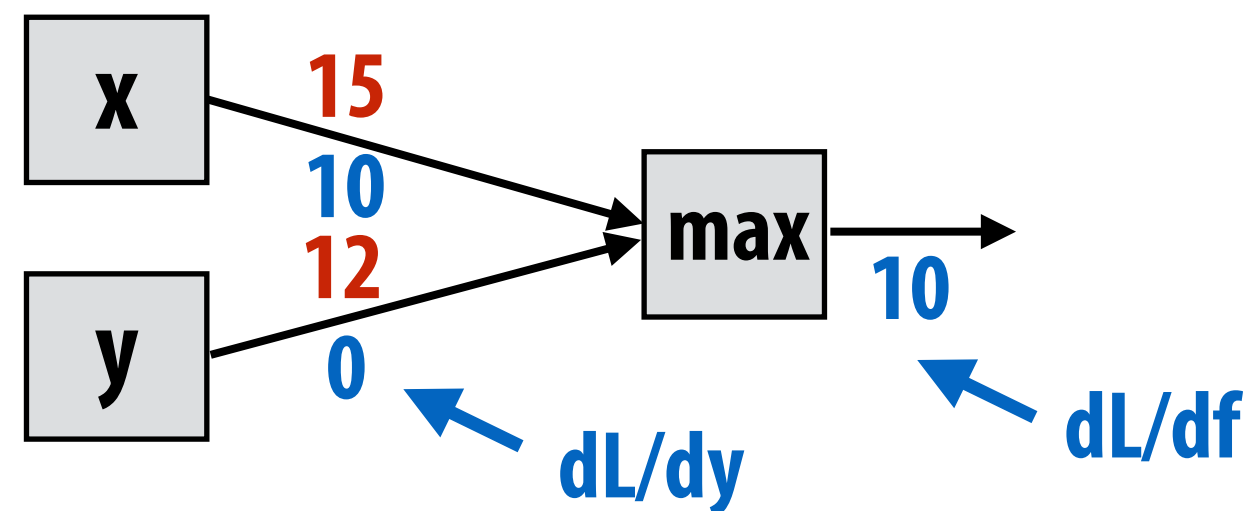
Back propagation as a reverse-mode autodiff technique

Red = output of node
Blue = $df/dnode$

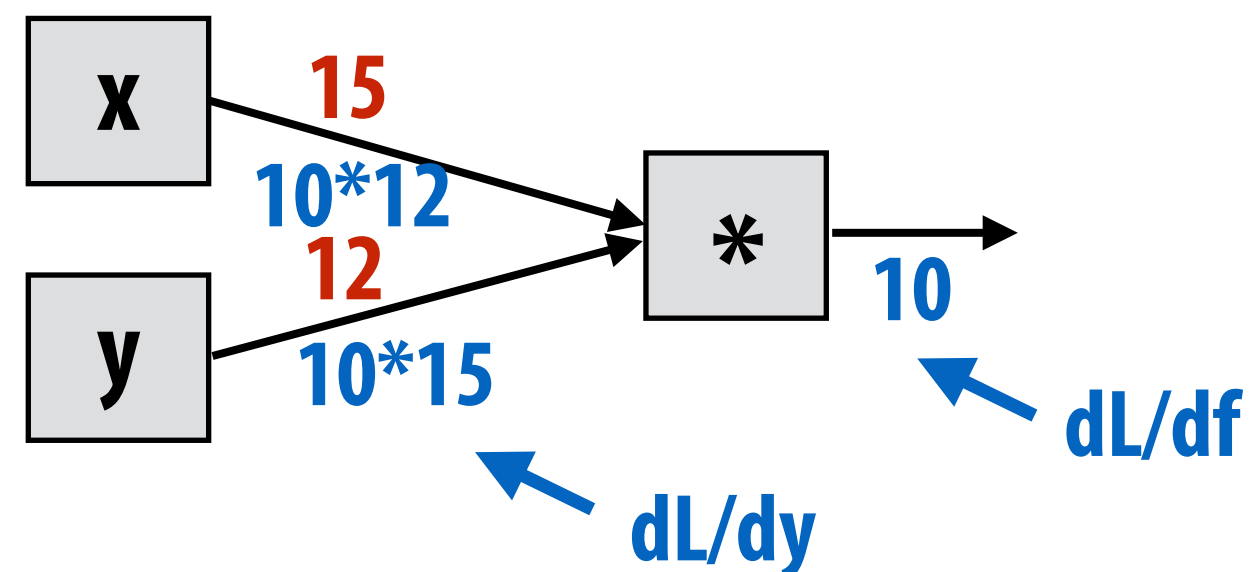
Recall: $\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$



$$f(x, y) = x + y \quad \frac{df}{dx} = 1, \frac{df}{dy} = 1$$



$$f(x, y) = \max(x, y) \quad \frac{df}{dx} = \begin{cases} 1, & \text{if } x > y \\ 0, & \text{otherwise} \end{cases}$$

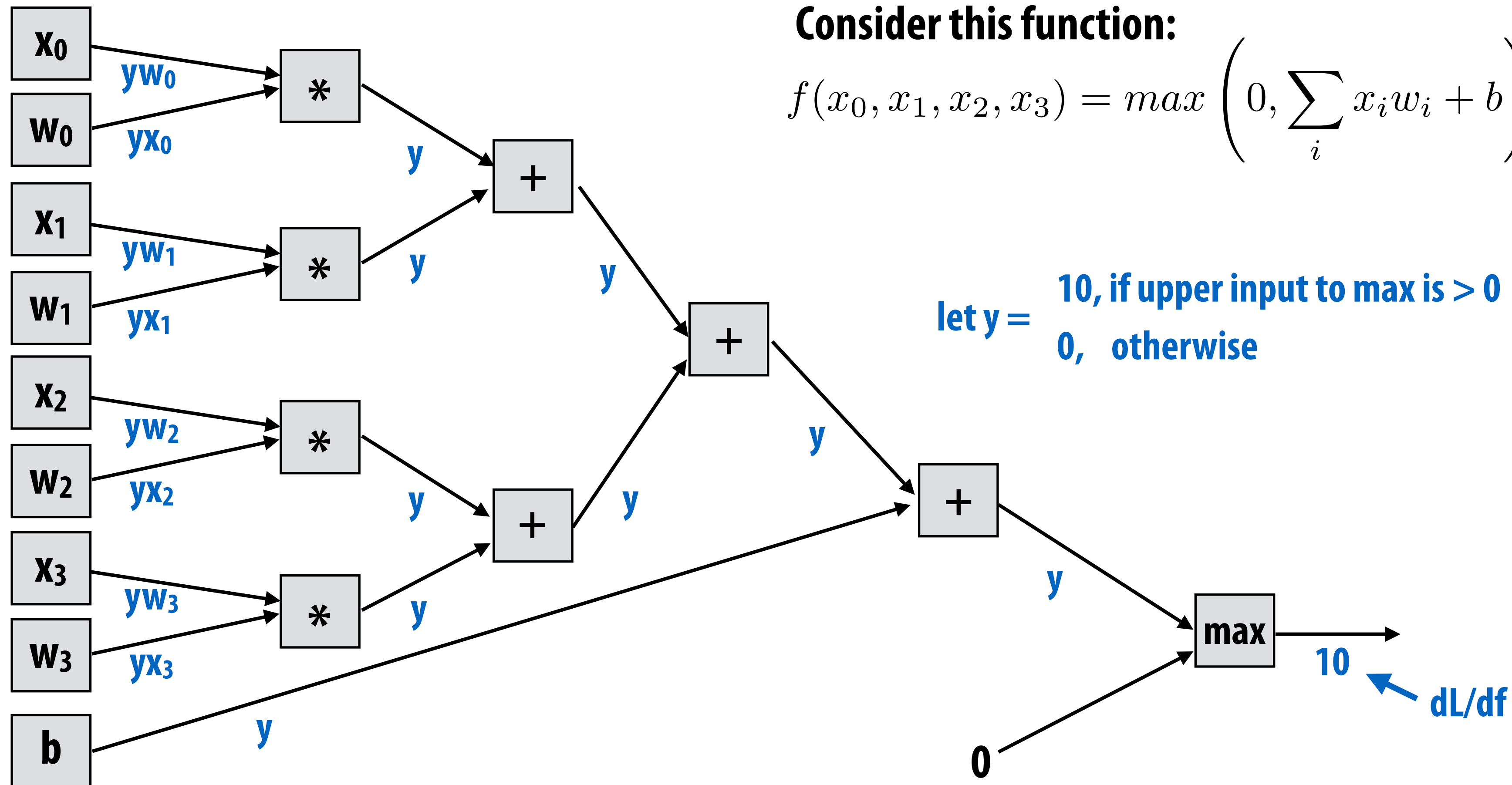


$$f(x, y) = xy \quad \frac{df}{dx} = y, \frac{df}{dy} = x$$

Back-propagation through a more complex function

Consider this function:

$$f(x_0, x_1, x_2, x_3) = \max \left(0, \sum_i x_i w_i + b \right)$$



Observe: output of prior node must be retained by the algorithm in order to compute gradients (with respect to weights) for this function during backprop.

Reverse-mode auto-diff in Slang

```
[Differentiable] // Auto-diff requires that functions are marked differentiable
float2 foo(float a, float b)
{
    return float2(a * b * b, a * a);
}
```

← Explicitly tell compiler differentiation support is desired

```
void main()
{
```

```
    DifferentialPair<float> dp_a = diffPair(
        1.0 // input 'a'
    ); // Calling diffPair without a derivative part initializes to 0.
```

← Value of parameter 'a'

```
    DifferentialPair<float> dp_b = diffPair(2.4);
```

← Value of parameter 'b'

```
    // Derivatives of scalar L w.r.t output.
```

```
    float2 dL_doutput = float2(1.0, 0.0);
```

← DL/df

```
    // bwd_diff to compute dL_da and dL_db
```

```
    // The derivative of the output is provided as an additional _input_ to the call
```

```
    // Derivatives w.r.t inputs are written into dp_a.d and dp_b.d
```

```
    //
```

```
    bwd_diff(foo)(dp_a, dp_b, dL_doutput);
```

← Call to generated bwd_diff function will compute gradients of foo's output with respect to a and b, and write them to dp_a.d and dp_b.d

```
    // Extract the derivatives of L w.r.t input
```

```
    float dL_da = dp_a.d;
```

```
    float dL_db = dp_b.d;
```

Custom derivative function overrides

- When derivative is difficult/expensive/impossible to derive mechanically, programmer can just specify the function to compute the derivatives directly.

```
void sin_bwd(inout DifferentialPair<float> dpx, float dresult)
{
    float x = dpx.p;

    // Write-back the derivative to each input (the primal part must be copied over as-is)
    dpx = DifferentialPair<float>(x, cos(x) * dresult);
}

[BackwardDerivative(sin_bwd)]
float sin(float x)
{
    // Calc sin(X) using Taylor series..
}
```

For a sophisticated example use case:

see “Differentiable Vector Graphics Rasterization for Editing and Learning” Li et al. 2020

Using interfaces to mix differentiable/non-differentiable code

■ Differentiable interfaces

```
struct MyType : IDifferentiable
{
    float x;
    float y;
};
```



```
struct MyType : IDifferentiable
{
    // Automatically inserted by Slang from the fact that
    // MyType has 2 floats which are both differentiable
    //
    typealias Differential = MyType;
    // ...
}
```

← **Compiler infers the type of the Differential
(which in this case is also a pair of floats)**

```
MyType obj;
obj.x = 1.f;

MyType.Differential d_obj;
// Differentiable fields will have a corresponding field in the diff type
d_obj.x = 1.f;
```

**And now its valid for the programmer to create
instances of either MyType, or MyType.Differential**

```
struct MyPartialDiffType : IDifferentiable
{
    // Automatically inserted by Slang based on which fields are differentiable.
    typealias MyPartialDiffType = syn_MyPartialDiffType_Differential;

    float x;
    uint y;
};
```

```
// Synthesized
struct syn_MyPartialDiffType_Differential
{
    // Only one field since 'y' does not conform to IDifferentiable
    float x;
};
```



**Compiler infers the type of the Differential
(Now just a single float)**

```
struct MyPartialDiffType : IDifferentiable
{
    // Automatically inserted by Slang based on which fields are differentiable.
    typealias MyPartialDiffType = syn_MyPartialDiffType_Differential;

    float x;
    uint y;
};

// Synthesized
struct syn_MyPartialDiffType_Differential
{
    // Only one field since 'y' does not conform to IDifferentiable
    float x;
};
```

Using interfaces to mix differentiable/non-differentiable code

- Programmer can tell compiler which function parameters to differentiate over

```
// Only differentiate this function with regard to `x`.  
float myFunc(no_diff float a, float x);
```



```
void back_prop(float a, inout DifferentialPair<float> x, float dResult);
```

Summary

- **Slang is designed to be extremely UN-OPINIONATED about how the programmer goes about authoring differentiable code**
- **It's a set of programming language constructs and a compiler implementation that handles the “busy work” of generating gradient code, and type checking code to ensure the programmer is not making statically-checkable mistakes**
 - **With reasonable “defaults” but almost everything is overridable**

Project presentations

Presentation slots are being set up on Tues and Wed

- **Tuesday:**

- **10:00-noon**
- **1:30-2:50**

- **Wednesday:**

- **3:00-4:00**

- **We hope these slots work for most teams, if you have a conflict, let's work it out case-by-case**

- **Location TBD**

- **Everyone is invited to come watch during any of these times**

Project handins

- **Final writeup is due Tuesday 5pm (no exceptions)**
 - **Guidelines/instructions online**
- **Oral presentation:**
 - **8-minute presentation per team. We will cut you off to keep time! :(**
 - **We expect all members of the team to have a speaking role**
- **Why presentations and not posters?**
 - **Kayvon believes communicating a technical idea clearly in a short amount of time is a far more important life skill than standing in from of a poster**
 - **Unfortunately, most engineers do this pretty poorly... let's work on it!**

Communicating like an architect

(Aka. project presentation tips)

A few course themes

■ *Thinking like a systems architect*

- What are inputs/outputs, constraints, and goals?
- What are the “services” the system should perform (what is hard for a user to do in a world without the system)
- Once you can establish answers to these questions, you can consider your solution options

■ Knowledge of applications and systems is necessary to choose efficient solutions

- Algorithms folks use their hammer to reduce cost of algorithms (smaller DNN models, new optimization hyperparameters, use different data structures, etc.)
- HW designers use their hammer (custom accelerators, new interconnects) to execute a given workload faster
- SW systems folks use their hammer (parallel/distributed computing, workload-specific scheduling, high level programming abstractions, etc.)

■ The best solutions architects pick the right mixture of hammers for the job

A few course themes

- **Knowledge of applications and systems is necessary to do meaningful evaluation**
 - **Is this the right workload to evaluate?**
 - **Pitfall: measuring speedup on a part of the workload that is not the most significant**
 - **Does dataset I'm using have the right data distribution?**
 - **Am I measuring cost in FLOPs, but increasing data movement?**
 - **Am I optimizing an algorithm that is not the right choice of algorithm for this problem?**

Why get good at communication?

- **Systems architects need to be some of the best communicators in an organization**
- **Must be able to:**
 - **Communicate with users to understand workloads and constraints**
 - **Communicate with the individual contributors/engineers to:**
 - **Understand emerging problems/constraints**
 - **Understand how to evolve/extending existing designs to enable new functionality**
 - **Must constantly be communicating to various parties why:**
 - **Their desired features are not being added**
 - **They must do the same work with 50% of the resources**
 - **Communicate to everyone a strategic vision for a system**
 - **Communicate to executives/management/funders how goals are being met.**

My motivation

- **I have found I give nearly the same feedback over and over to students making talks**
 - **It is not profound feedback, it is just application of a simple set of techniques and principles that are consistently useful when making talks**
- **I am hoping these slides serve as a useful checklist you can refer to vet your own project presentation talks before giving them next Tuesday**
 - **Don't worry: I still make these mistakes all the time when creating first drafts of talks**

Who painted this painting?



Salvador Dali (age 22)

My point: learn the basic principles before you consciously choose to break them



Put yourself in your audience's shoes

Tip 1

This is a major challenge for most technical speakers. (including professors)

(Tip: recite a sentence out loud to yourself. * Do you really expect someone who has not been working with you everyday on the project to understand what you just said?)

*** I'm not kidding. Say it out loud. I find hearing myself say something out loud makes it easier to parse it from an audience's perspective.**

Consider your audience

- **Everyone in the audience knows about course readings/topics**
 - Terminology/concepts we all know about need not be defined (just say “remember we talked about X”)
- **Most of the audience knows little-to-nothing about the specific application domain or problem you are trying to solve**
 - Application-specific terminology should be defined or avoided
 - What they can get their head around is inputs/outputs and goals/constraints
- **Everyone wants to know the “most interesting” thing that you found out or accomplished (your job is to define most interesting for them)**

Tip 2

**A good principle for any talk (or paper):
“Every sentence matters”**

What are you trying to say?

What technical story are you trying to tell?

What is point you are trying to make?

You might be trying to define a hypothesis.

Or define goals

Or establish inputs/outputs

**Or show data that suggests you
were successful.**

Is what you just said making that point? (If not, remove it)

If you can't justify how it will help the listener understand the point, take it out.

If it's not likely to be clear, take it out.

Pick a focus

- In this class, different projects should stress different results
- Some projects may wish to show a flashy demo and describe how it works (proof by “it works”)
- Other projects may wish to show a sequence of graphs (path of progressive optimization) and describe the optimization that took system from performance A to B to C
- Other projects may wish to clearly contrast parallel CPU vs. parallel GPU performance for a workload

Your job is not to explain what you did, but to explain what you think we should know.

**And really the most important thing we want to know is (1) what was your goal? and
(2) what's the evidence you have that you were successful?**


Ignoring every sentence matters

Never ever, ever, ever do this!

Outline
<ul style="list-style-type: none">• Introduction• Related Work• Proposed System Architecture<ul style="list-style-type: none">❖ Basic design decision❖ Dedicated hardware for T&I❖ Reconfigurable processor for RGS• Results and Analysis• Conclusion

Bad example 2


■ Who is the audience for this? (how does this benefit them?)



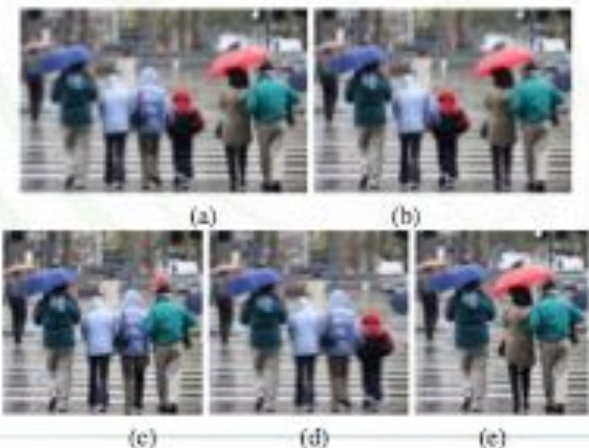
Discrete Methods

- Remove/Add Pixels
- Related work*:
 - Seam Carving [SIGGRAPH 07]
 - Improved Seam Carving [SIGGRAPH 08]
 - Shift-Map Image Editing [ICCV 09]

* non-exhaustive



From: Avidan et al., Seam Carving for Content-Aware Image Resizing




From: Pritch et al., Shift-Map Image Editing

This type of related work section says little more than “others have worked in this area before”.


- I suspect your audience assumes this is the case.
- Every sentence matters: if it doesn't provide value, take it out (or replace it with comments that do provide value)

- Experts?
 - They likely know these papers exist. These slides don't tell them what about these papers is most relevant to this talk
- Non-experts?
 - They won't learn the related work from these two slides




Continuous methods


- Find continuous transformation
- Warp/deformation grid
- Related work*
 - Non-homogenous warping, ICCV 07
 - Streaming video, SIGGRAPH 09
 - Shrinkability Maps for Content-Aware Video Resizing, PG 09
 - Robust Image Retargeting via Axis-Aligned Deformation, EG 12



Generated with the Streaming Video approach [KRA09]



original saliency map



retargeting to 200% width using axis-aligned deformations

From: Robust Image Retargeting via Axis-Aligned Deformation

Tip 3

The audience prefers not to think (much)

The audience has a finite supply of mental effort

- **The audience does not want to burn mental effort about things you know and can just tell them.**
 - They want to be led by hand through the major steps of your story
 - They do not want to interpret any of your figures or graphs, they want to be directly told how to interpret them (e.g., what to look for in a graph).
 - They want to be told about your key assumptions
- **The audience does want to spend their energy thinking about:**
 - Potential problems with what you did (did you consider all edge cases? Is your evaluation methodology sound? Is this a good platform for this workload?)
 - Implications of your approach to other things
 - Connections to their own project or interests

Tip 4

Set up the problem.

**Establish inputs, outputs, and constraints
(goals and assumptions)**

Basics of problem setup

- **What is the computation performed (or system built)?**
 - What are the inputs? What are the outputs?
- **Why does this problem stand to benefit from optimization?**
 - “Real-time performance could be achieved”
 - “Researchers could run many more trials, changing how science is done”
 - “It is 90% of the execution time in this particular system”
- **Why is it hard? (What made your project interesting? What should we reward you for?)**
 - What turned out to be the hardest part of the problem?
 - Optimization projects: e.g., overcoming SIMD divergence, increasing arithmetic intensity
 - Applications projects: coming up with an algorithm to estimate depth, dealing with motion blur
 - Abstraction/API projects: deciding between stateless or stateful interfaces

Example: 3D rendering problem

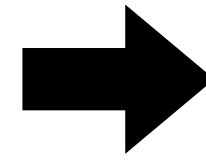
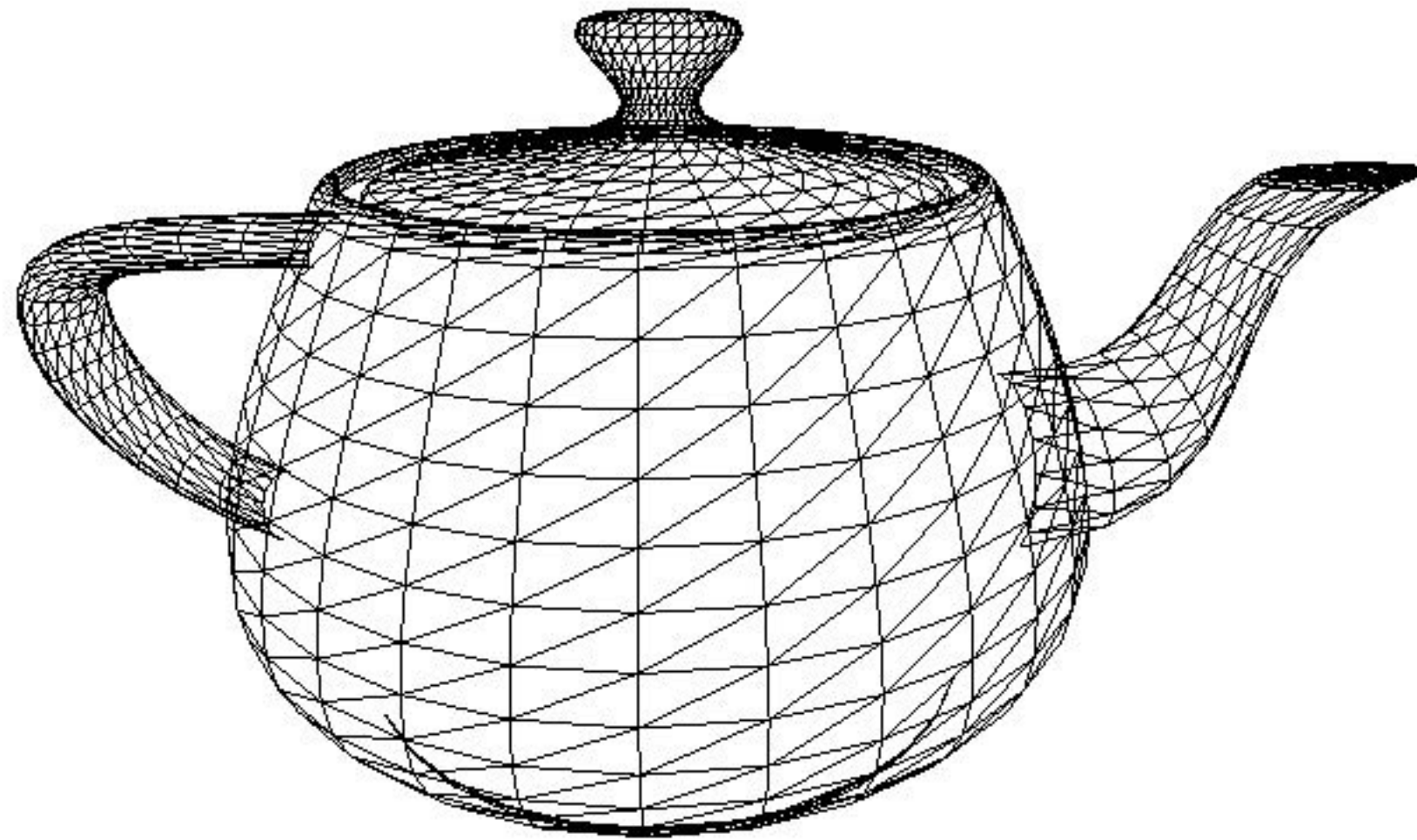


Image credit: Henrik Wann Jensen

Input: description of a scene:

3D surface geometry (e.g., triangle mesh)
surface materials, lights, camera, etc.

Output: image of the scene

Simple definition of rendering task: computing how each triangle in 3D mesh contributes to appearance of each pixel in the image?

Why is knowing the goals and constraints important?

Your contribution is typically a system or algorithm that meets the stated goals under the stated constraints.

Understanding whether a solution is “good” requires having this problem context.

Tip 5

How to describe a system

How to describe a system

■ Start with the nouns (the key boxes in a diagram)

- Major components (processors, memories, interconnects, etc.)
- Major entities (particles, neighbor lists, pixels, pixel tiles, features, etc.)
- What is state in the system?

■ Then describe the verbs

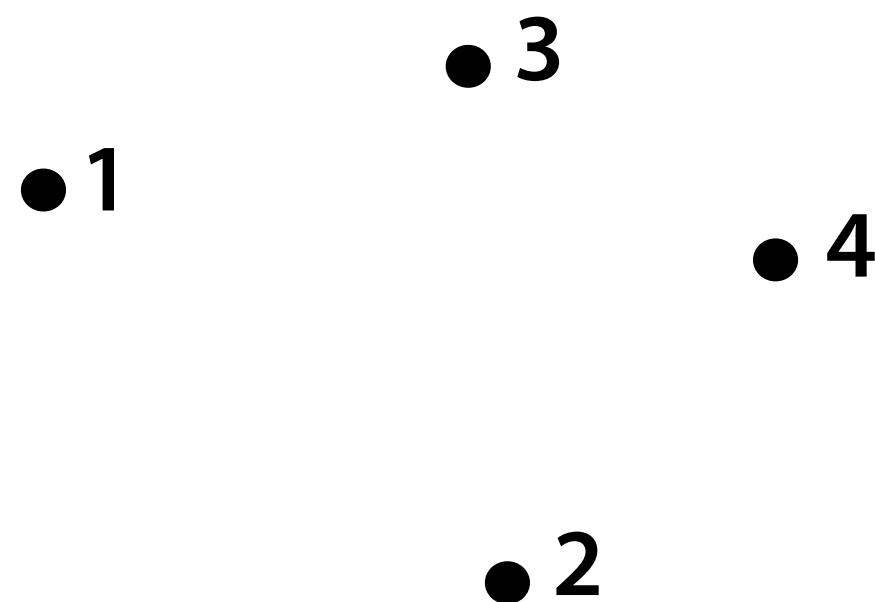
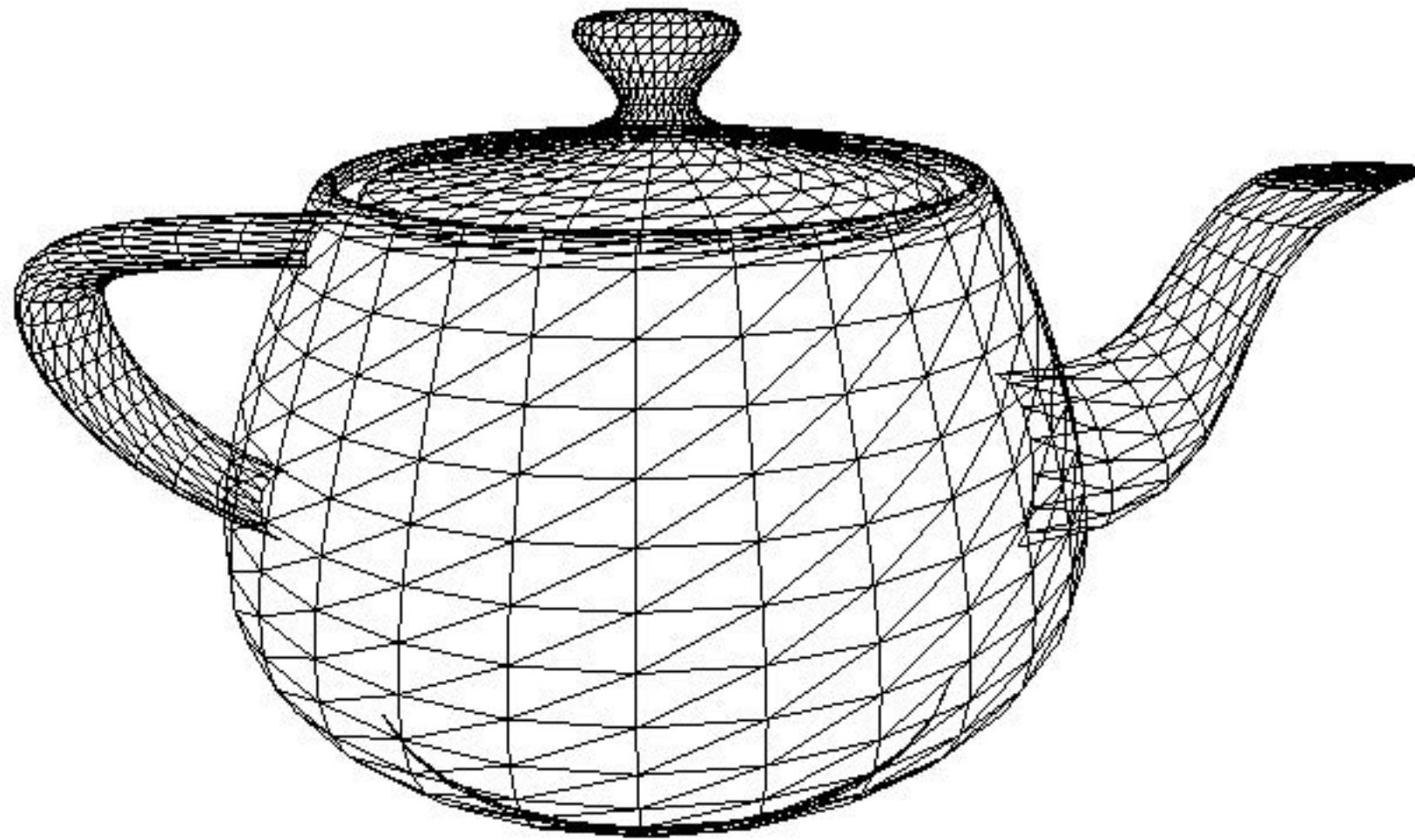
- Operations that can be performed on the state (update particle positions, compute gradient of pixels, traverse graph, etc.)
- Operations produce, consume, or transform entities

Tip: how to explain “a system”

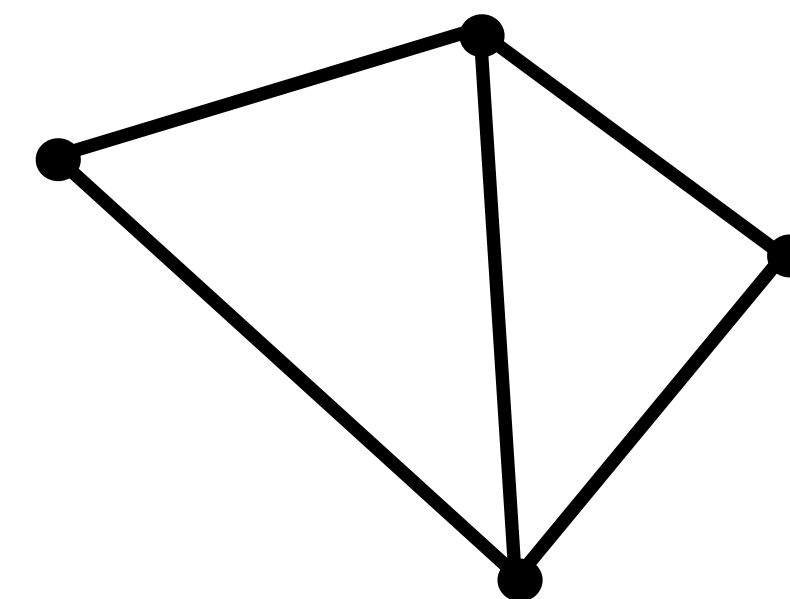
- Step 1: describe the things (key entities) that are manipulated
 - The nouns

Example: real-time graphics primitives (entities)

Represent surface as a 3D triangle mesh

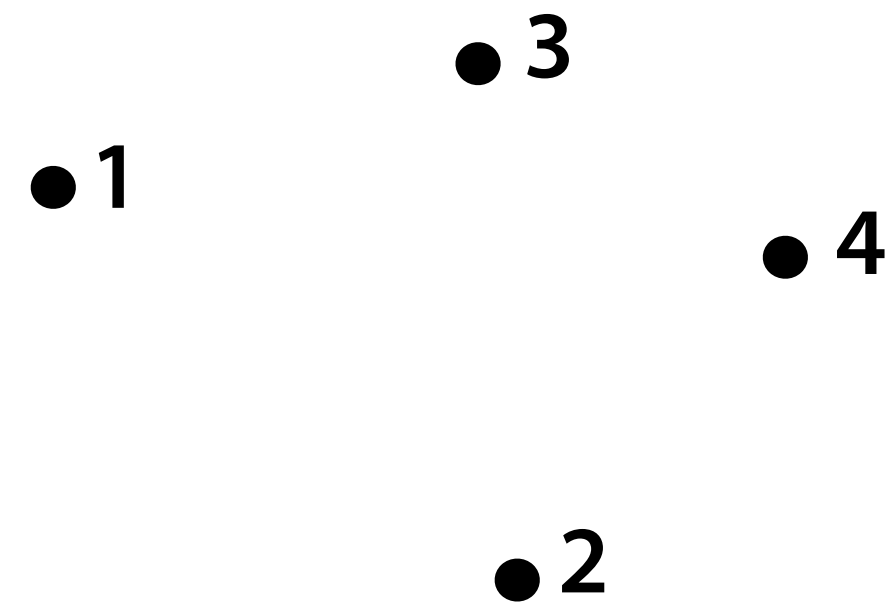


Vertices
(points in space)

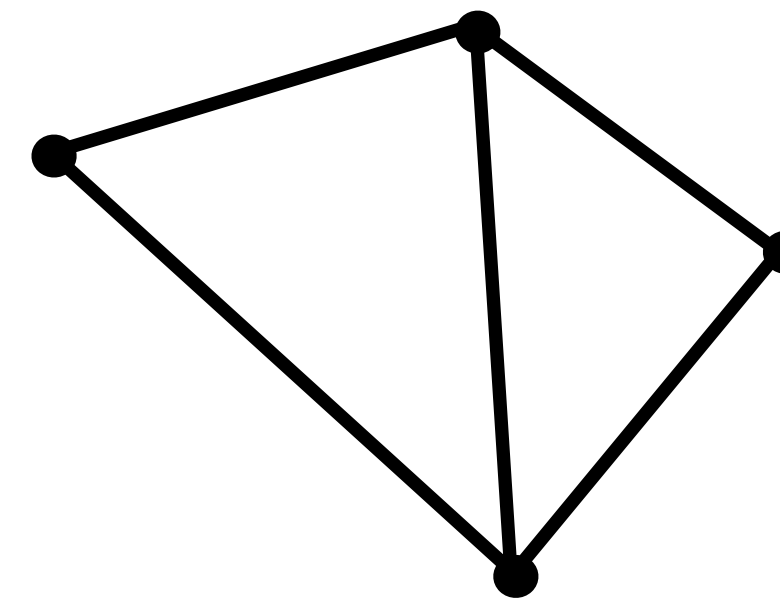


Primitives
(e.g., triangles, points, lines)

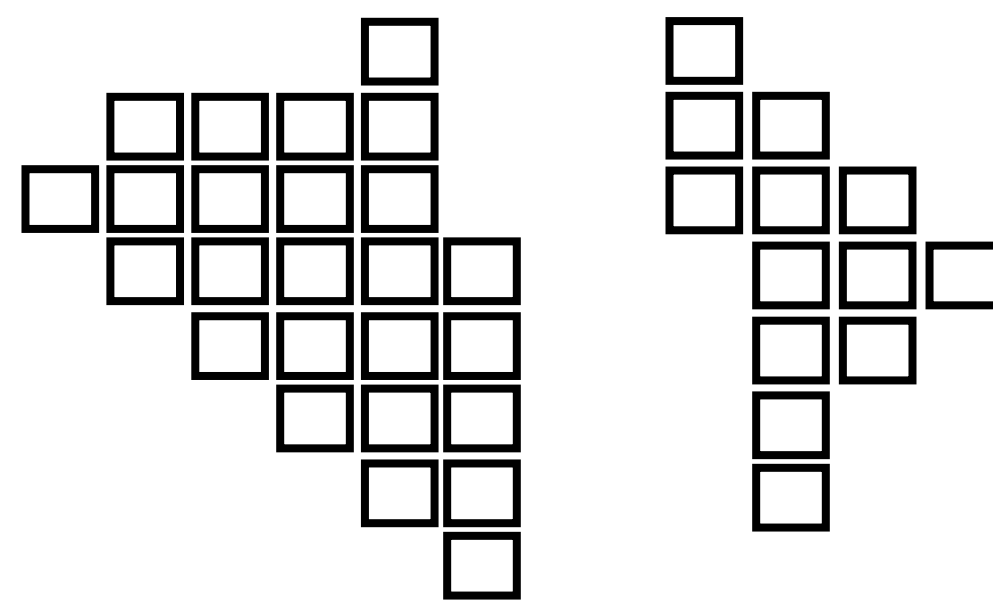
Real-time graphics primitives (entities)



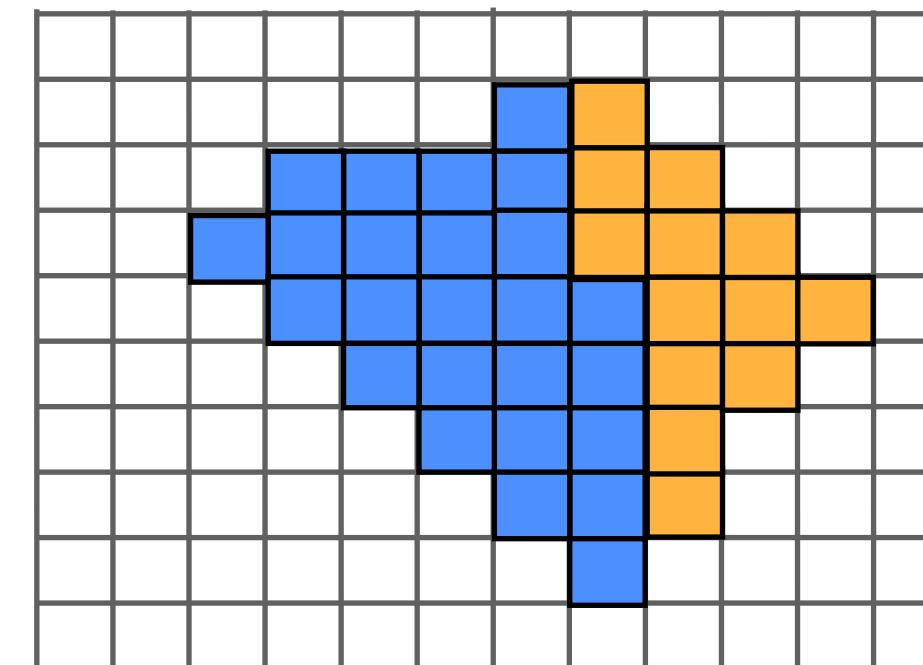
Vertices
(points in space)



Primitives
(e.g., triangles, points, lines)



Fragments



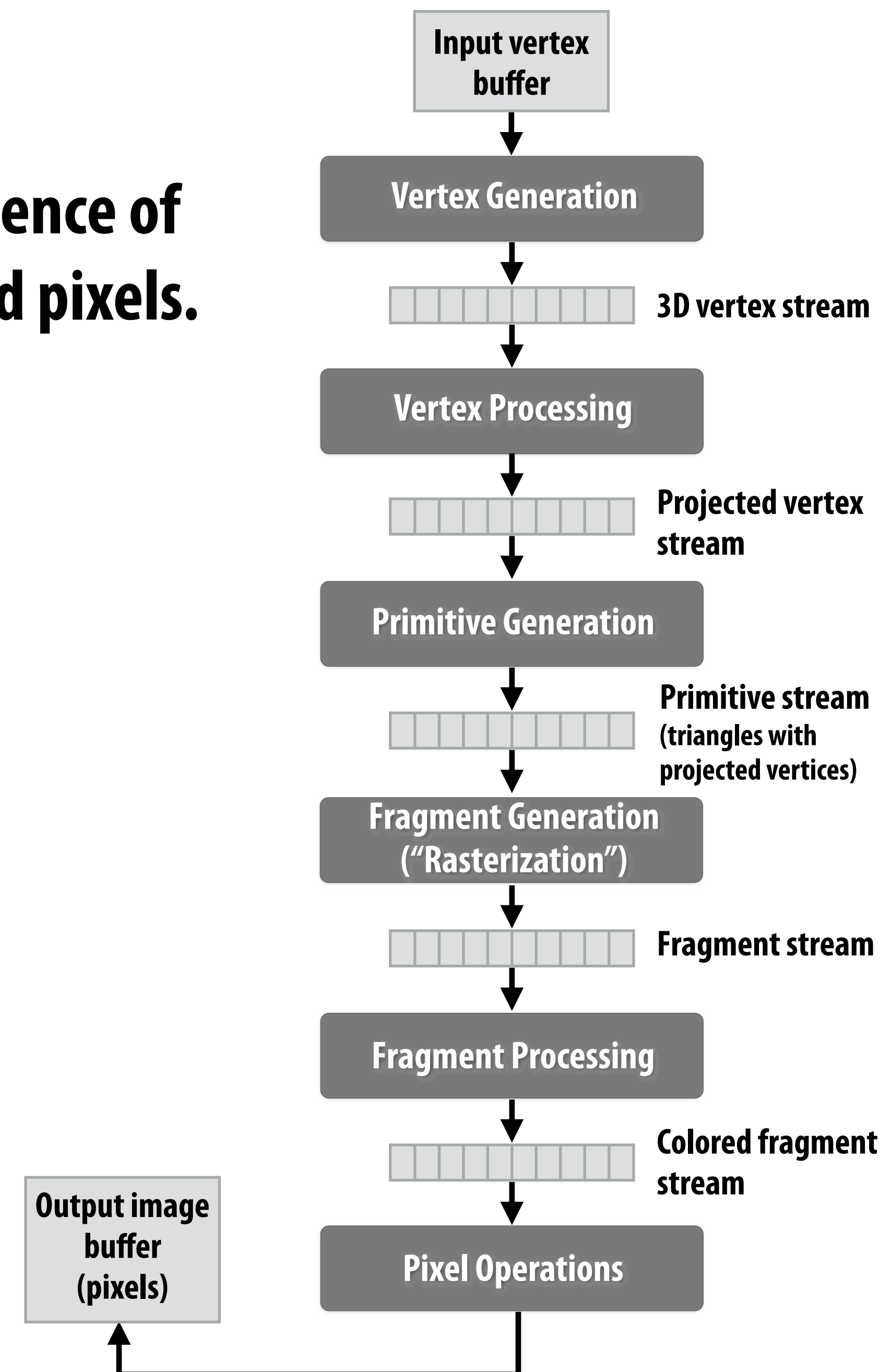
Pixels (in an image)

How to explain “a system”

- **Step 1: describe the things (key entities) that are manipulated**
 - **The nouns**
- **Step 2: describe the operations the system performs on these entities**
 - **The verbs**

Real-time graphics pipeline

- Abstracts process of rendering a picture as a sequence of operations on vertices, primitives, fragments, and pixels.



Tip 6

Surprises* are almost always bad:

Say where you are going and why you must go there before you say what you did.

*** I am referring to surprises in talk narrative and/or exposition. A surprising result is great.**

Give the why before the what

■ Why provides the listener context for...

- Compartmentalizing: assessing how hard they should pay attention (is this a critical idea, or just an implementation detail?). Especially useful if they are getting lost.
- Understanding how parts of the talk relate (“Why is the speaker now introducing a new optimization framework?”)

■ In the algorithm description:

- “We need to first establish some terminology”
- “Even given X , the problem we still haven’t solved is...”
- “Now that we have defined a cost metric we need a method to minimize it...”

■ In the results/evaluation:

- Speaker: “Key questions to ask about our approach are...”
- Audience: “Thanks! I agree, those are good questions. Let’s see what the results say!”

Two key questions:

- How much does SRDH improve traversal cost when perfect information about shadow rays is present?
- How does the benefit of the SRDH decrease as less shadow ray information is known a priori? (Is a practical implementation possible?)

Big surprises in a narrative are a bad sign

- Ideally, you want the audience to always be able to anticipate* what you are about to say
 - This means: your story is so clear it's obvious!
 - It also means the talk is really easy to present without notes or text on slides (it just flows)
- If you are practicing your talk, and you keep forgetting what's coming on the next slide (that is, you can't anticipate it)...
 - This means: you probably need to restructure your talk because a clear narrative is not there.
 - It's not even obvious to you! Ouch!

* Credit to Abhinav Gupta for suggesting the term anticipation, and for the example on this slide

Tip 7

Show, don't tell

**It's much easier to communicate with
figures/images than text**

(And it saves the speaker a lot of work explaining... you can just describe the picture)

Example:

- In a recent project, we asked the question... given enough video of tennis matches of a professional athlete, could we come up with an algorithm for turning all this input video into a controllable video game character?

Compare the description above to the following sequence...

Here's an example of that source video

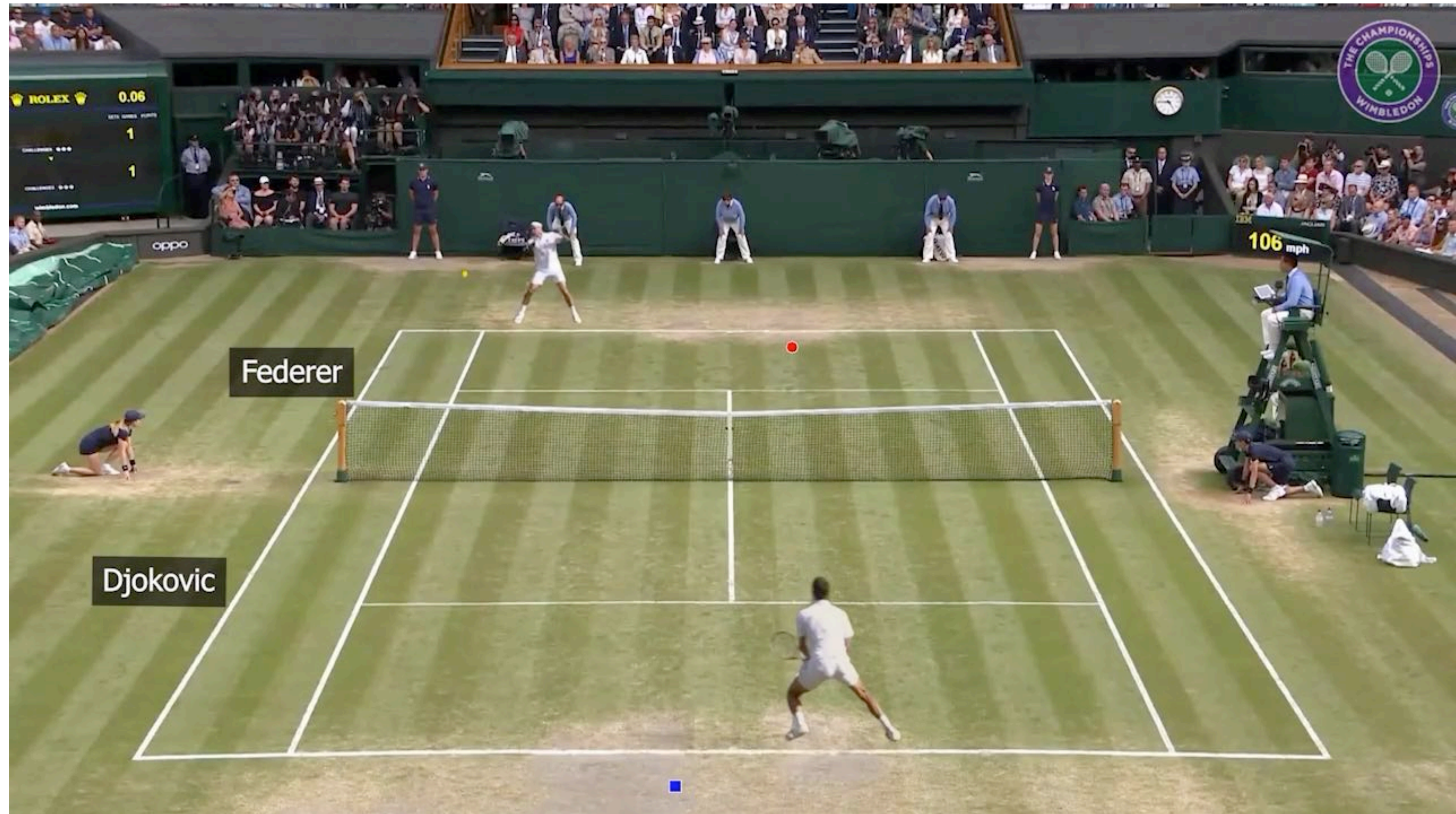


The best way to describe the input data is just show it! ("This is what the input looks like!")

And there's a lot of it out there!



And here's an example of controllable output



The best way to describe the output we seek is just show the result of the system!
("We click to specify a target ball location, and the player hits the incoming ball back to the red dot")

Another example:



The problem (lighting differences)

After the fix



Another example: we recently created a renderer that achieved high frame rates by rendering many views of the scene at the same time



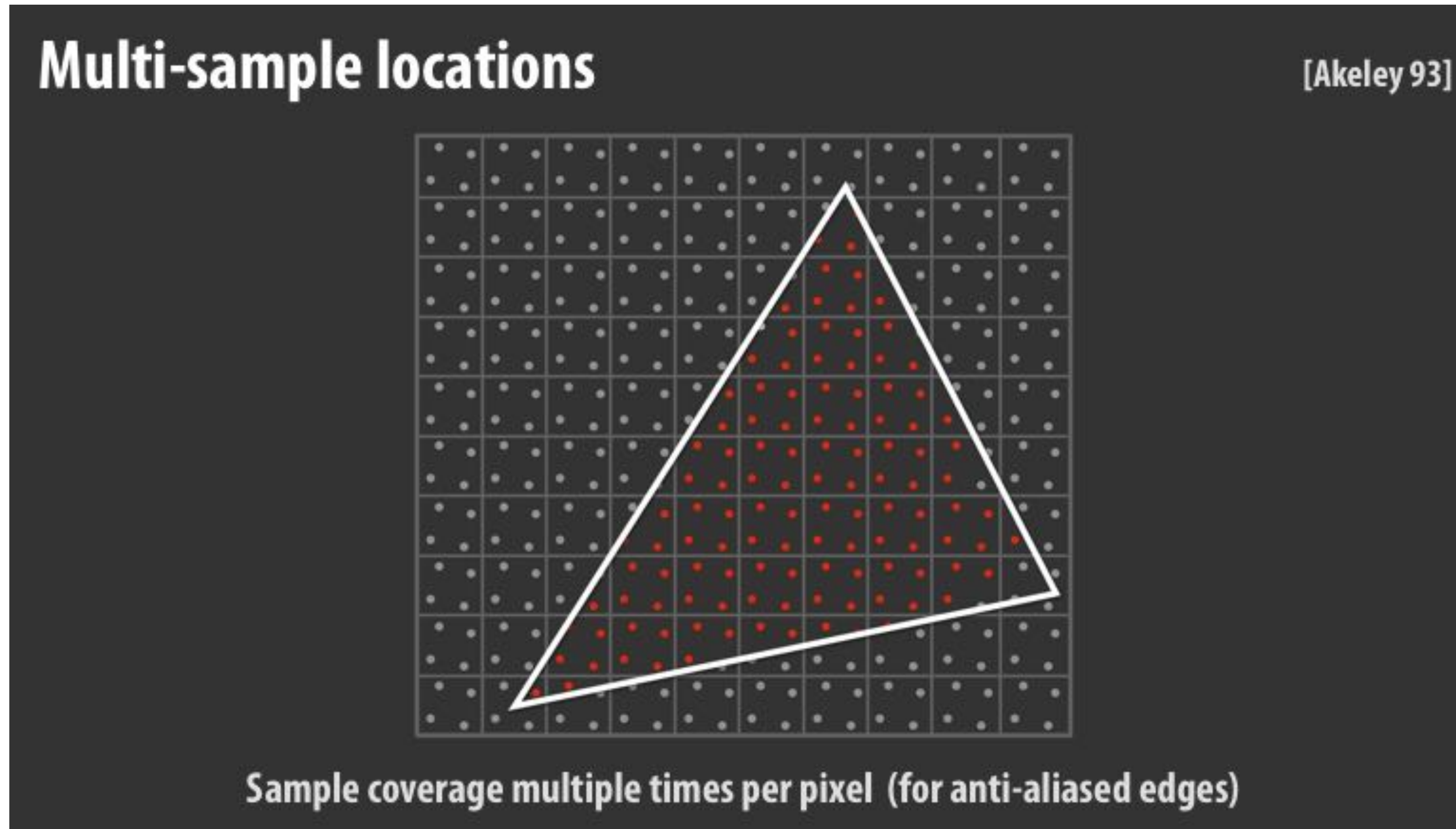
Tip 8

***Always, always, always
explain any figure or graph***

(remember, the audience does not want to think about things you can tell them)

Explain every figure

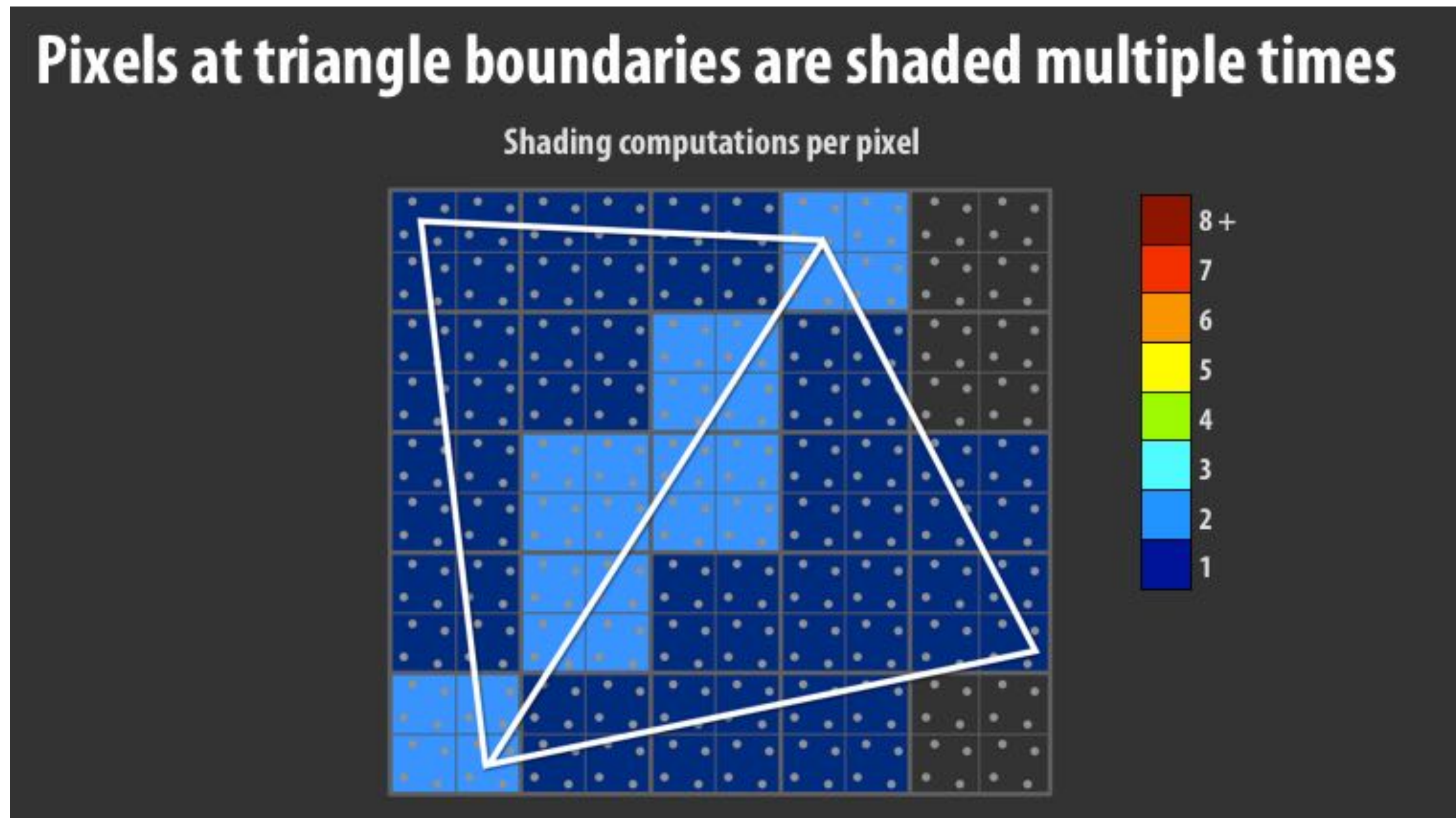
- Explain every visual element in the figure (never make the audience decode a figure)
- Refer to highlight colors explicitly (explain why the visual element is highlighted)



Example voice over: “Here I’m showing you a pixel grid, a projected triangle, and the location of four sample points at each pixel. Sample points falling within the triangle are colored red.”

Explain every figure

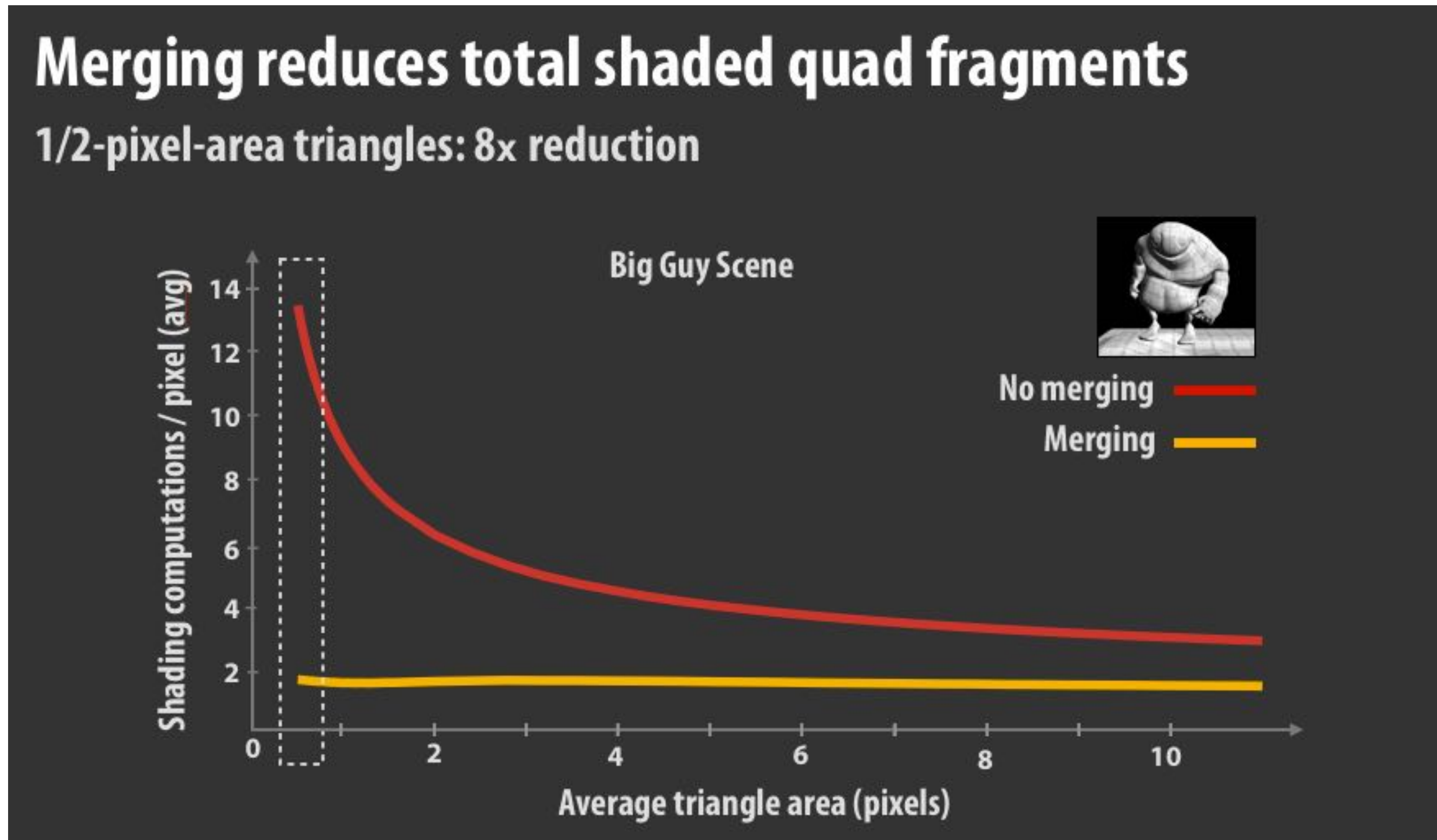
- Lead the listener through the key points of the figure
- Useful phrase: “As you can see...”
 - It’s like verbal eye contact. It keeps the listener engaged and makes the listener happy... “Oh yeah, I can see that! I am following this talk!”



Example voice over: “Now I’m showing you two adjacent triangles, and I’m coloring pixels according to the number of shading computations that occur at each pixel as a result of rendering these two triangles. As you can see from the light blue region, pixels near the boundary of the two triangles get shaded twice.

Explain every results graph

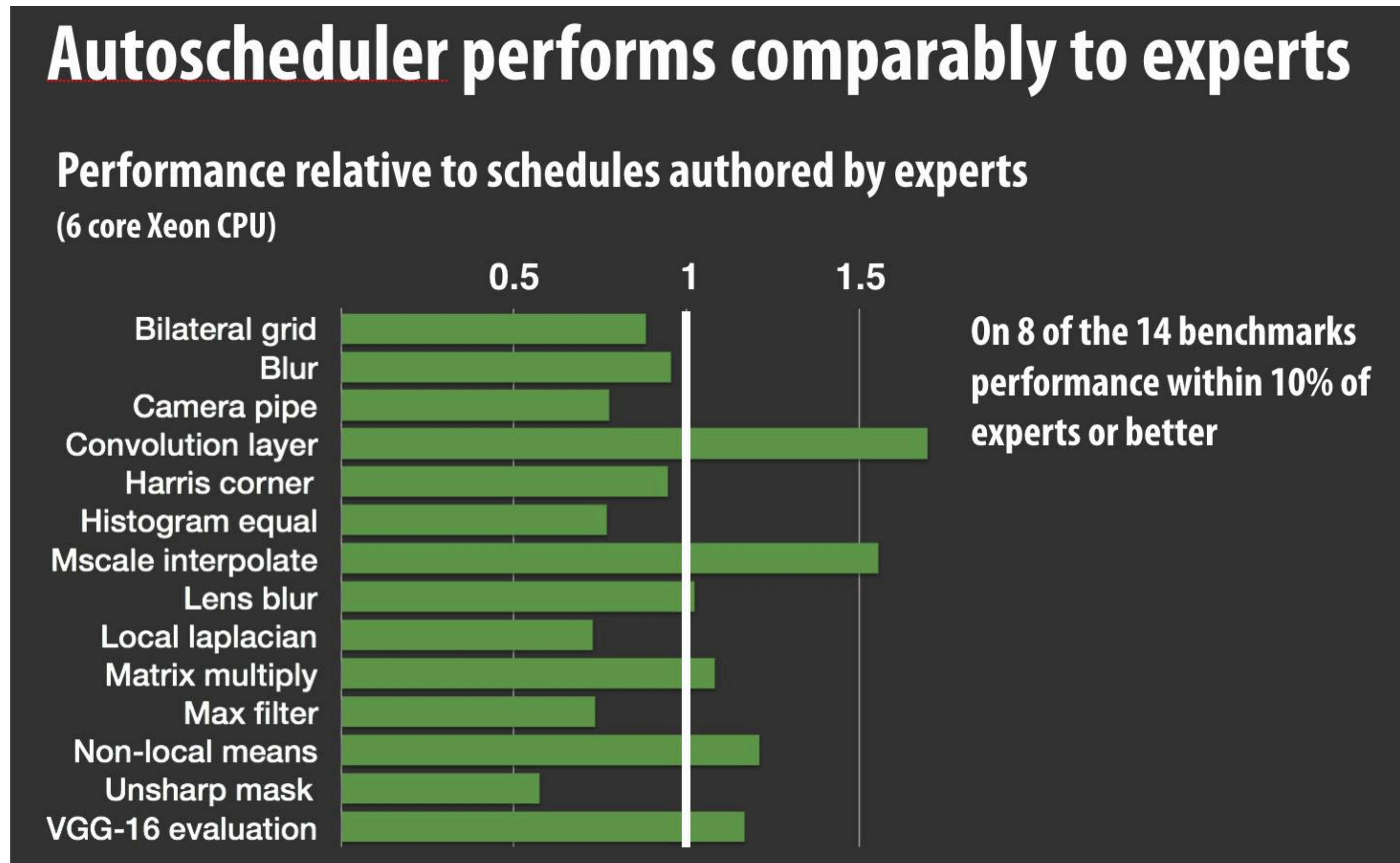
- May start with a general intro of what the graph will address (so audience “anticipates” the result)
- Then describe the axes (and your axes better have labels!)
- Then describe the one point that you wish to make with this results slide (more on this later!)



Example voice over: “Our first questions were about performance: how much did the algorithm reduce the number of the shading computations? And we found out that the answer is a lot. This figure plots the number of shading computations per pixel when rendering different tessellations of the big guy scene. X-axis gives triangle size. If you look at the left side of the graph, which corresponds to a high-resolution micropolygon mesh, you can see that merging, shown by yellow line, shades over eight times less than the convention pipeline.

Explain every results graph

- May start with a general intro of what the graph will address.
- Then describe the axes (your axes better have labels!)
- Then describe the one point that you wish to make with this results slide (more on this later!)



Example voice over: “Our first question was about performance: how fast is the auto scheduler compared to experts? And we found out that it’s quite good. This figure plots the performance of the autoscheduler compared to that of expert code. So expert code is 1. Faster code is to the right. As you can see, the auto scheduler is within 10% of the performance of the experts in many cases, and always within a factor of 2.

Tip 9

In the results section:

One point per slide!

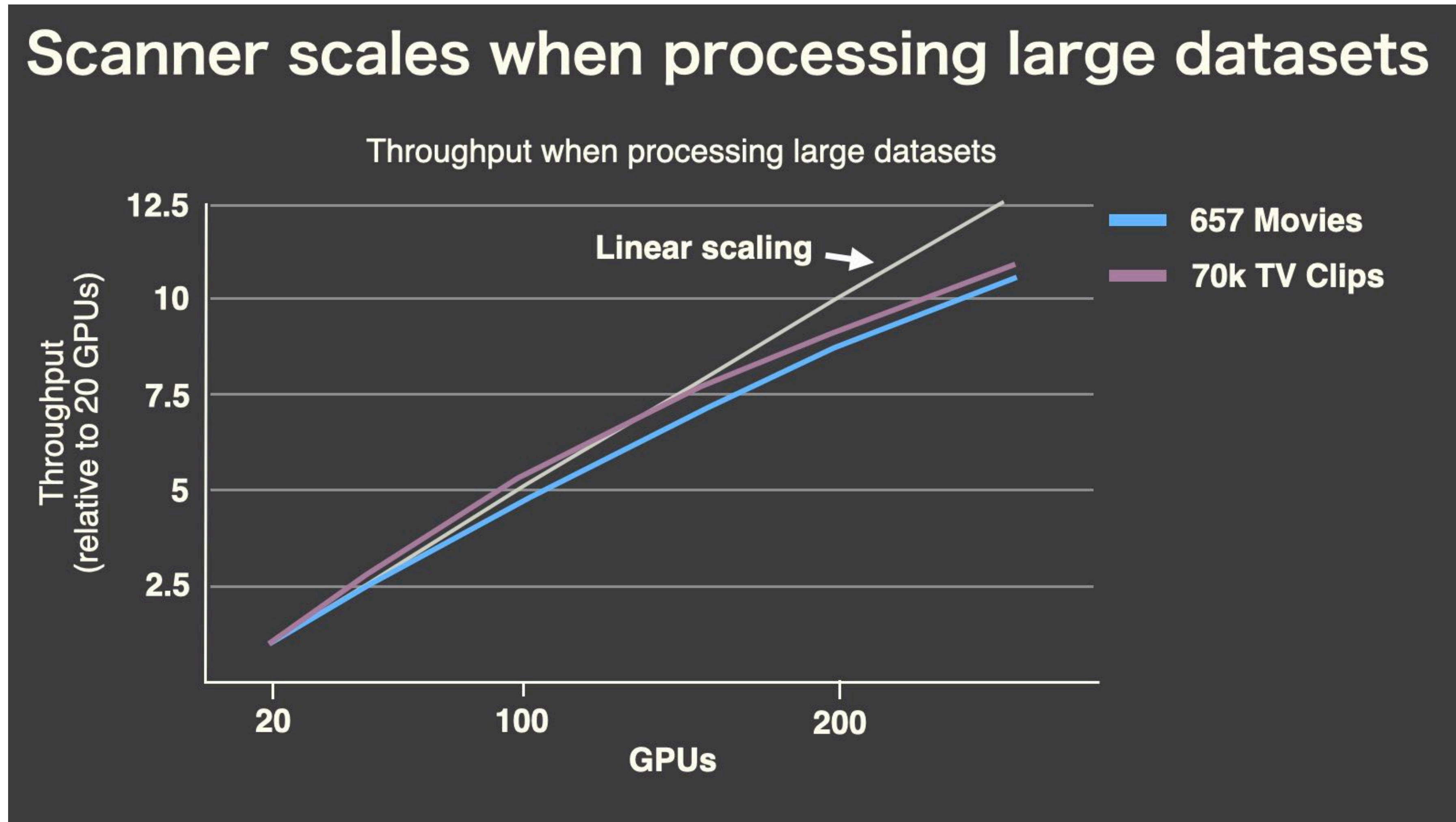
One point per slide!

One point per slide!

(and the point is the title of the slide!!!)

■ **Make the point of the graph the slide's title:**

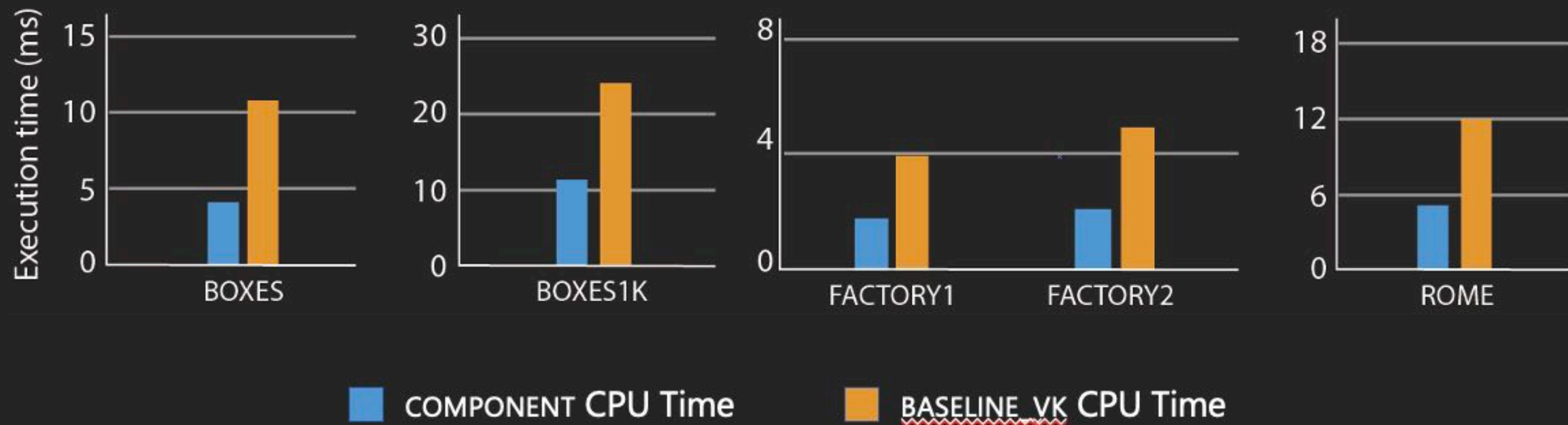
- **It provides audience context for interpreting the graph (“Let me see if I can verify that point in the graph to check my understanding”)**
- **Another example of the “audience prefers not to think” principle**



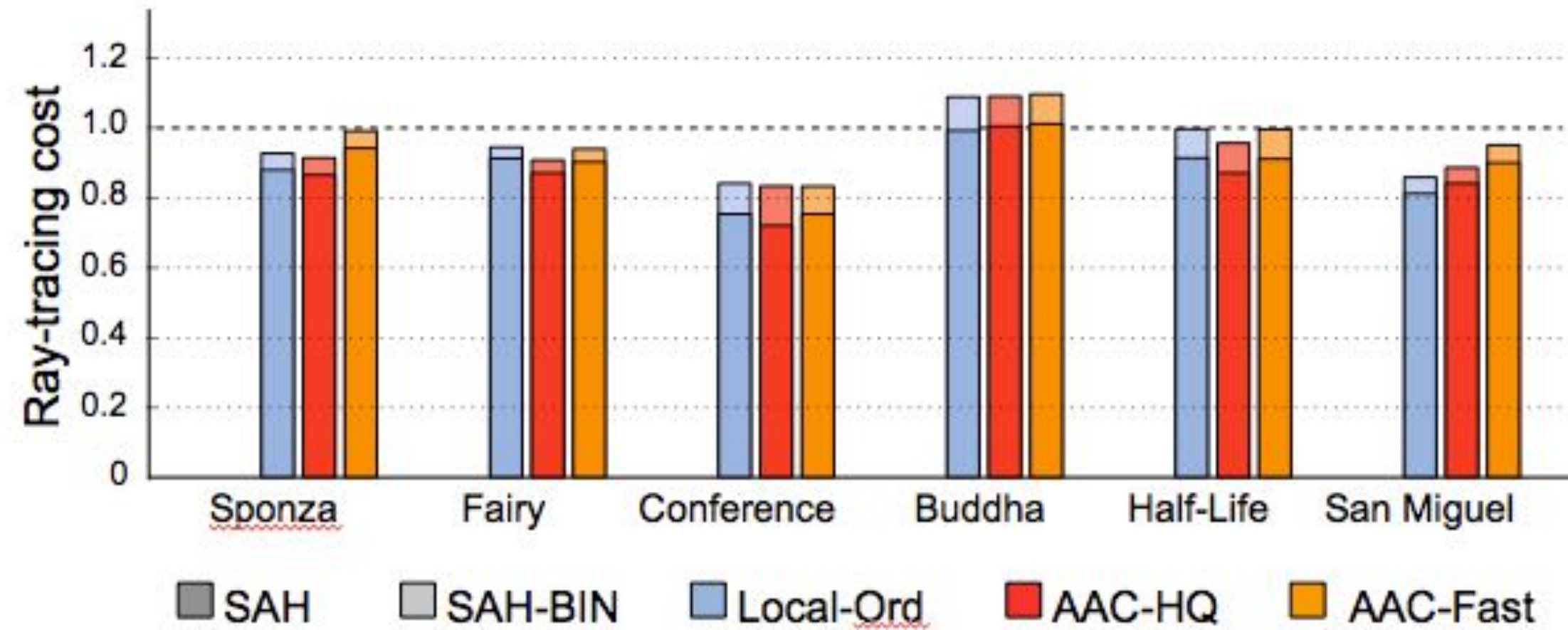
Another example:

The COMPONENTS renderer uses 2x less CPU time than BASELINE_VK

CPU Performance Comparison (single core)



AAC-Fast produces BVHs with equal or lower cost than the **full sweep build** in all cases except Buddha.

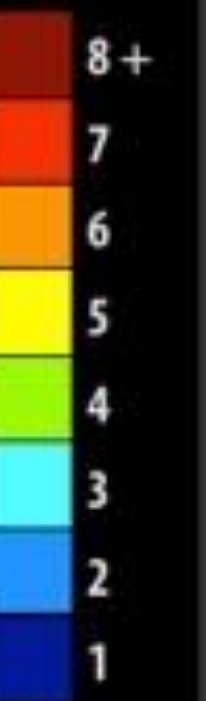


More examples

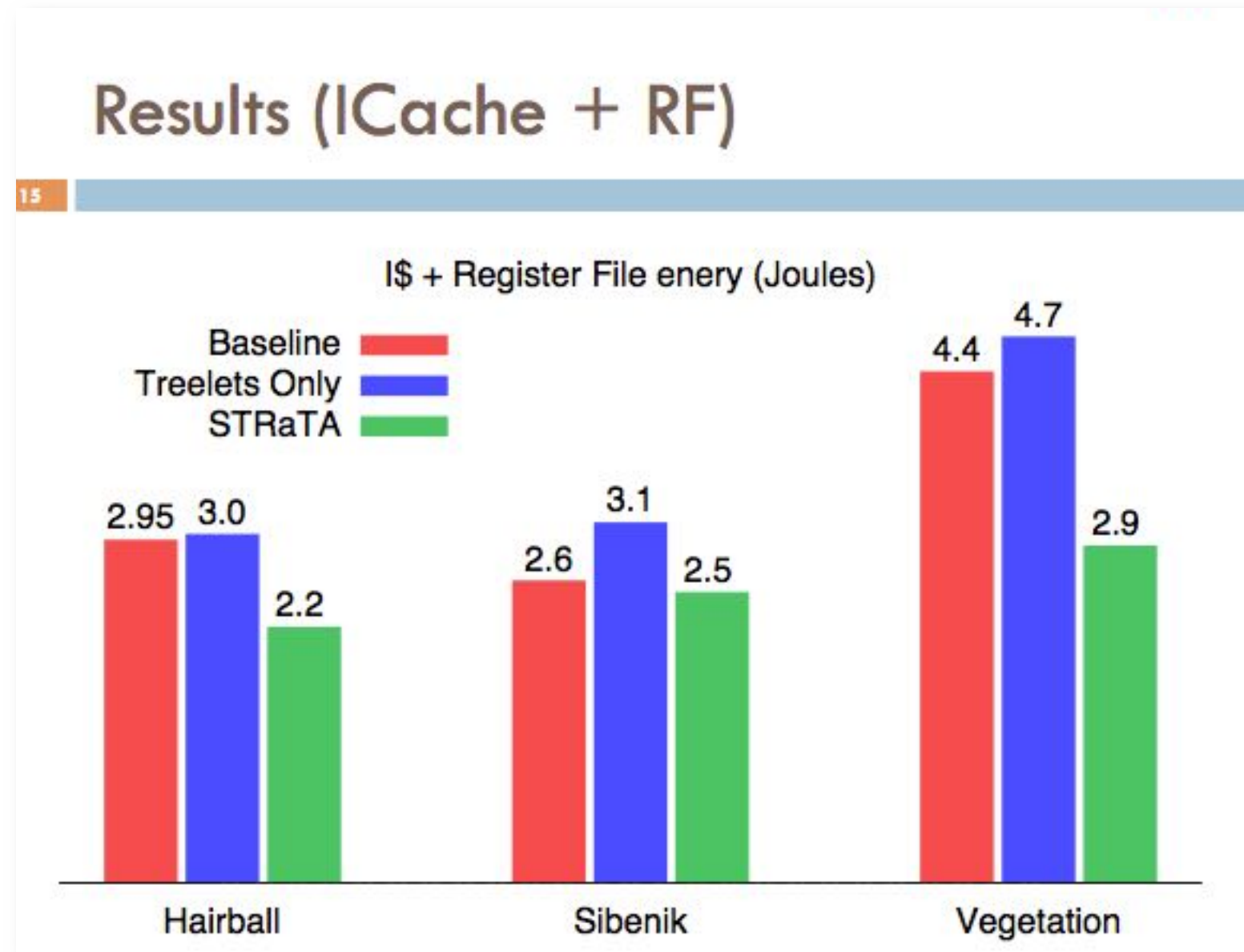
Extra shading occurs at merging window boundaries



1/2 pixel area triangles



Bad examples of results slides



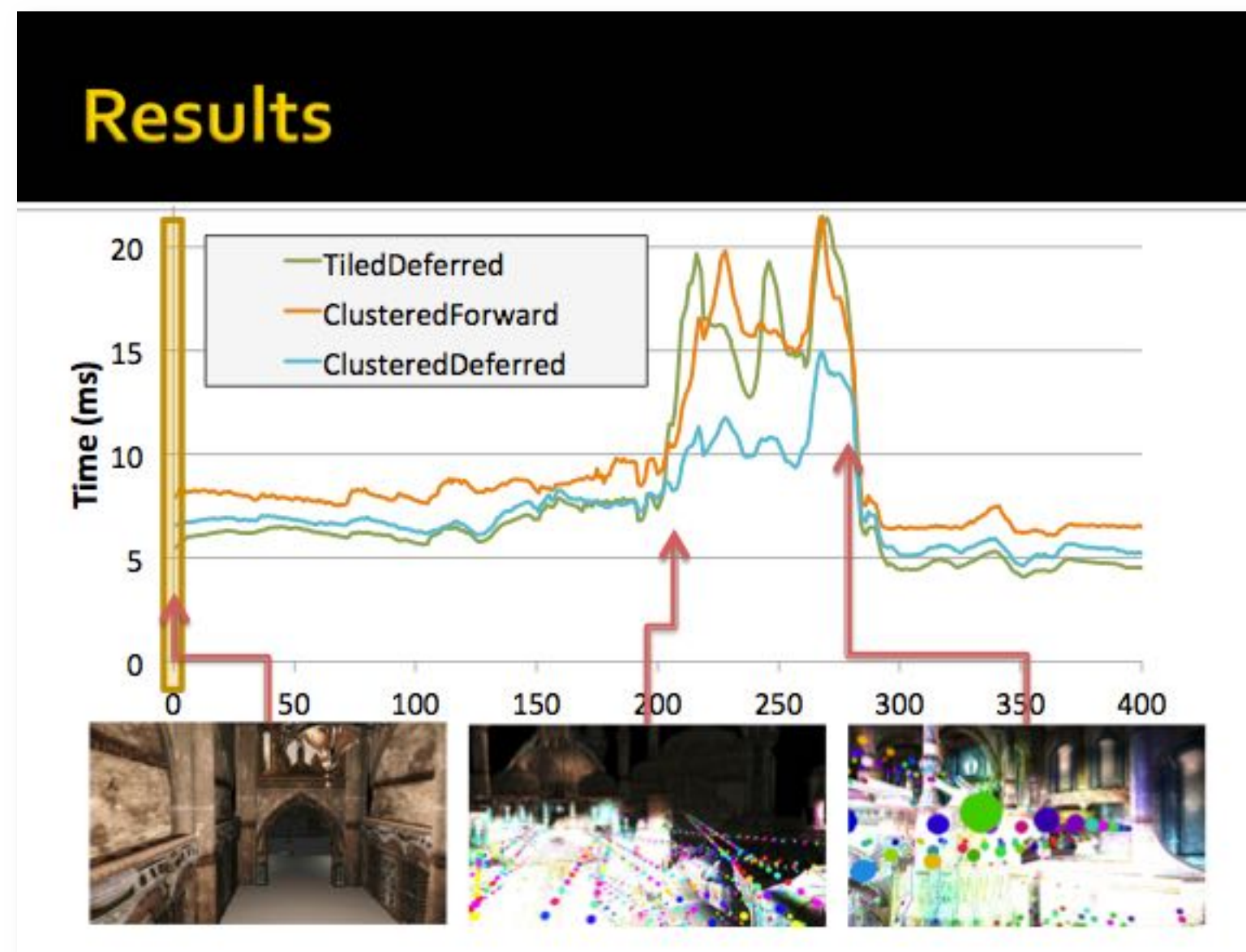
- Notice how you (as an audience member) are working hard to interpret the trends in these graphs
 - You are asking: what do these results say?
 - What am I supposed to be concluding?
- The audience just wants to be told what to look for!
 - They are reading the graphs to verify the main point, not determine the main point.

Simulation Results : RGS

RGS Performance

- ❖ 147-198 Mray/sec
- ❖ Texture cache concerns : Mip-mapping & Compression

Test scene	Ray type	Cache hit rate (%)		Bandwidth (GB/s)	Performance (Mrays/sec)
		Texture	Data		
Sibenik (80K tri.)	Primary	-	96.76	0.5	182.11
	FSR	-	91.24	1.9	172.25
Fairy (179K tri.)	Primary	93.25	96.87	0.8	175.66
	FSR	81.49	94.91	1.9	147.45
Ferrari (210K tri.)	Primary	86.12	98.09	0.6	183.28
	FSR	75.95	95.71	2.0	163.67
Conference (282K tri.)	Primary	-	98.44	0.2	198.32
	FSR	-	95.72	0.8	158.79



Tip 10

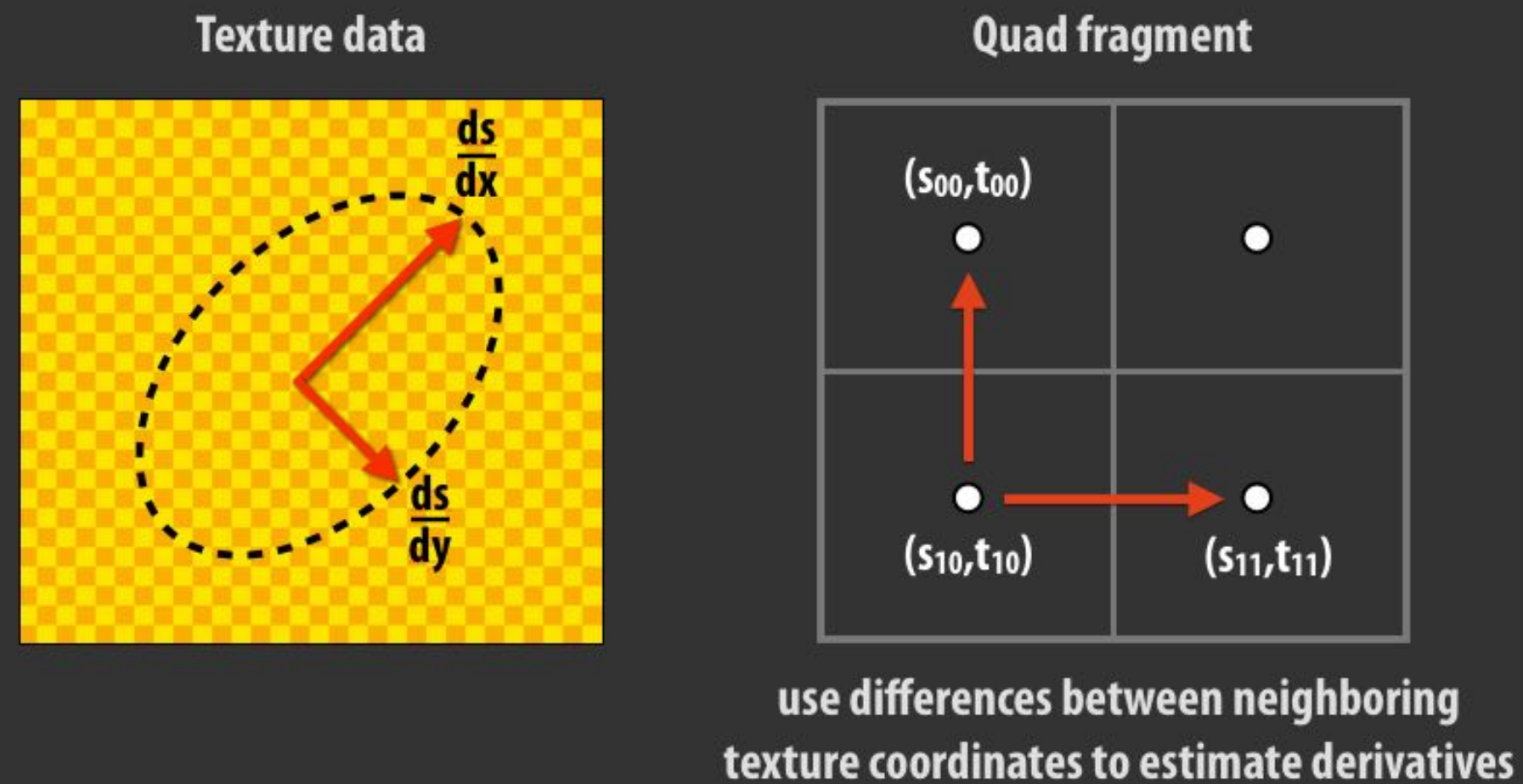
Titles matter

If you read the titles of your talk all the way through, it should be a great summary of the talk.

(basically, this is “one-point-per-slide” for the whole talk)

Examples of descriptive slide titles

GPUs shade quad fragments (2x2 pixel blocks)



Greedy SRDH build optimizes over partitions and traversal policies

SAH:

```
forall(partitions in set-of-partitions)
  ...evaluate SAH and pick min...
```

SRDH:

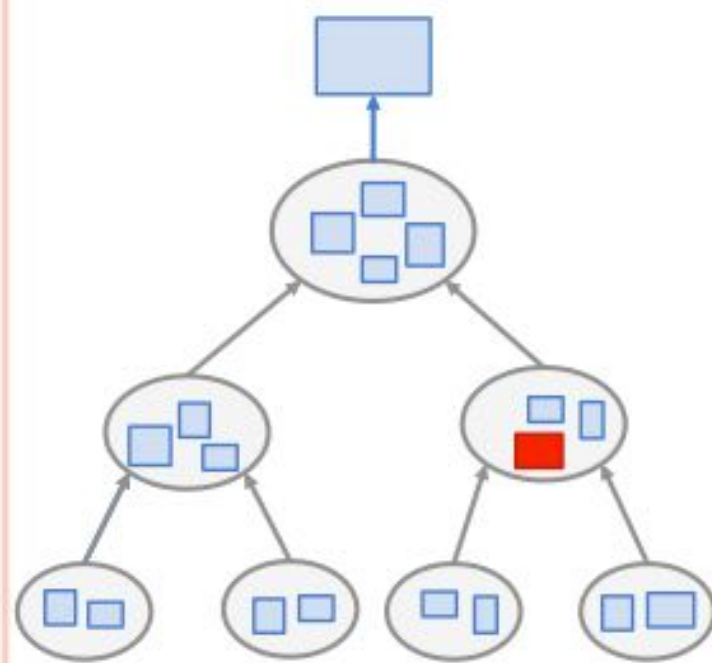
```
forall(partitions in set-of-partitions)
  forall(traversalKernels in set-of-kernels)
    ...evaluate SRDH and pick min...
```

$$\text{SRDH}(R, L, \kappa, r) = (1 - \kappa(r)H(L, r))|R| + (1 - \kappa(r)H(R, r))|L|$$

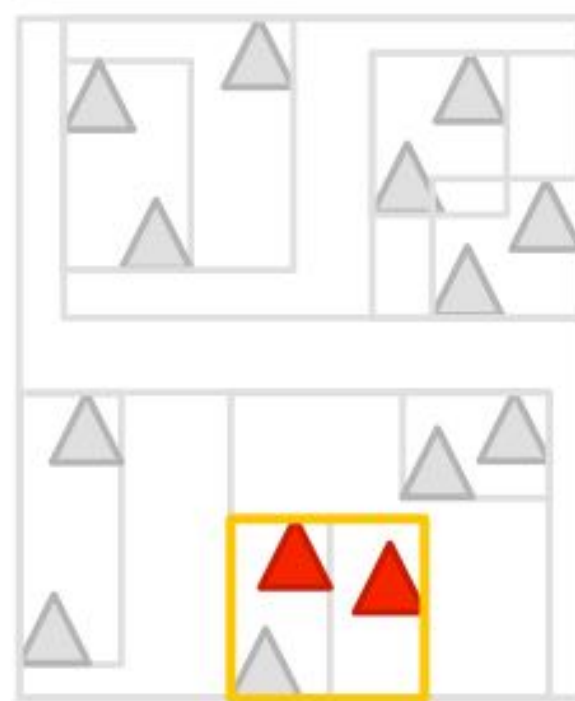
51

AAC IS AN APPROXIMATION TO THE TRUE AGGLOMERATIVE CLUSTERING SOLUTION.

Computation graph:



Primitive partitioning:



The reason for meaningful slide titles is convenience and clarity for the audience

“Why is the speaker telling me this again?”

(Recall “why before what”)

Read your slide titles in thumbnail view

Do they make all the points of the story you are trying to tell?

Reducing Shading on GPUs using Quad-Fragment Merging

Kayvon Fatahallian
Solomon Boalos
James Hegarty
Stanford University

Kurt Akeley
NVIDIA Research

Pat Hanrahan
Stanford University

William R. Mark
Intel Labs

Henry Morton
NVIDIA

1

High-resolution meshes are appearing in games

Low detail

2

High-resolution meshes are appearing in games

3

PROBLEM

Current GPUs shade small triangles inefficiently

4

Multi-sample locations

(Asterix 10)

5

Shading sample locations

(Asterix 10)

6

Surface derivatives are needed for texture filtering

Texture data

7

GPUs shade quad fragments (2x2 pixel blocks)

Texture data

Quad fragment

are differences between neighboring texture coordinates to estimate derivatives

8

Shaded quad fragments

9

Final pixel values

10

Pixels at triangle boundaries are shaded multiple times

Shading computations per pixel

11

Pixels at triangle boundaries are shaded multiple times

Shading computations per pixel

12

Pixels at triangle boundaries are shaded multiple times

Shading computations per pixel

13

Small triangles result in extra shading

Unmerged small triangles

Merged small triangles

7 pixel area triangles

14

Goal:

Shade high-resolution meshes (not individual triangles) approximately once per pixel

Approach:

Evolve GPU's quad-fragment shading system (Provide smooth evolution from status quo)

Not address alternatives later in talk: deferred shading, Ray-traced merged GPU shading

15

QUAD-FRAGMENT MERGING

16

GPU pipeline (with tessellation)

Mesh → Tessellation → Vert → Quad → Raster → Pixel

17

Rasterized quad-fragment

18

Rasterized quad-fragment

Shading inputs (interpolated from triangle vertices)

multi-sample coverage mask

19

Rasterized quad fragments

20

GPU pipeline: triangle connectivity is known

Mesh → Vert → Quad → Raster → Pixel

Triangle connectivity is known

quad-ops

21

Pipeline with quad-fragment merging

22

Pipeline with quad-fragment merging

23

Two key merging operations

1. Identifying when quad fragments can be merged

2. Constructing a merged quad fragment

24

Merging quad fragments

Mesh triangles

Rasterized quad fragments

Merged quad fragments

Step 1: aggregate coverage

25

Merging quad fragments

Mesh triangles

Rasterized quad fragments

Merged quad fragment

Step 2: sample shading inputs

Shading sample location

quad mask

26

Merging quad fragments

Mesh triangles

Rasterized quad fragments

Merged quad fragment

Step 2: sample shading inputs

Shading sample location

quad mask

27

Two key merging operations

1. Identifying when quad fragments can be merged

2. Constructing a merged quad fragment

28

Challenge

Avoiding merges that introduce visual artifacts

29

Example: surface with a silhouette

Triangle mask

Final pixels

anti-aliased silhouette

30

Naive merging results in aliasing

Triangle mask

Final pixels

aliasing result

31

Avoid merging across discontinuities

Mesh triangles

Merged quad fragments

Triangles 1, 2 (front-facing)

Triangles 3, 4 (back-facing)

Triangles 1, 3 (front-facing)

Triangles 2, 4 (front-facing)

Mask (top-left)

32

Conditions required to merge quad fragments

1. Same screen location

2. Same sidedness (triangles front-facing or back-facing)

3. Source triangles are adjacent in the mesh

33

High-frequency geometric detail may cause aliasing

Our merging rules are designed for real-time performance

Limit shading costs

Geometry should be pre-filtered to avoid aliasing

34

Implementation: the cost of merging is low

Merging operations are cheap

Testing merging rules requires only bitwise operations

Merge buffer is small

32 quad-fragment merge buffers is very effective

Expectation: quad-fragment merging can be encapsulated in fixed-function hardware

35

EVALUATION

36

Tip 11

Practice the presentation

Practice the presentation

- **Given the time constraints, you'll need to be smooth to say everything you want to say**
- **To be smooth you'll have to practice**
- **Rehearse your presentation several times the night before (in front of a partner or friend)**
 - **It's only a short presentation, so a couple of practice runs are possible in a small amount of time**