

Lecture 4:

Visibility

(coverage and occlusion using rasterization and the Z-buffer)

Visual Computing Systems
CMU 15-869, Fall 2014

Visibility

- **Very imprecise definition: computing what scene geometry is visible within each screen pixel**
 - What scene geometry projects into a screen pixel? (screen coverage)
 - Which geometry is actually visible from the camera at that pixel? (occlusion)

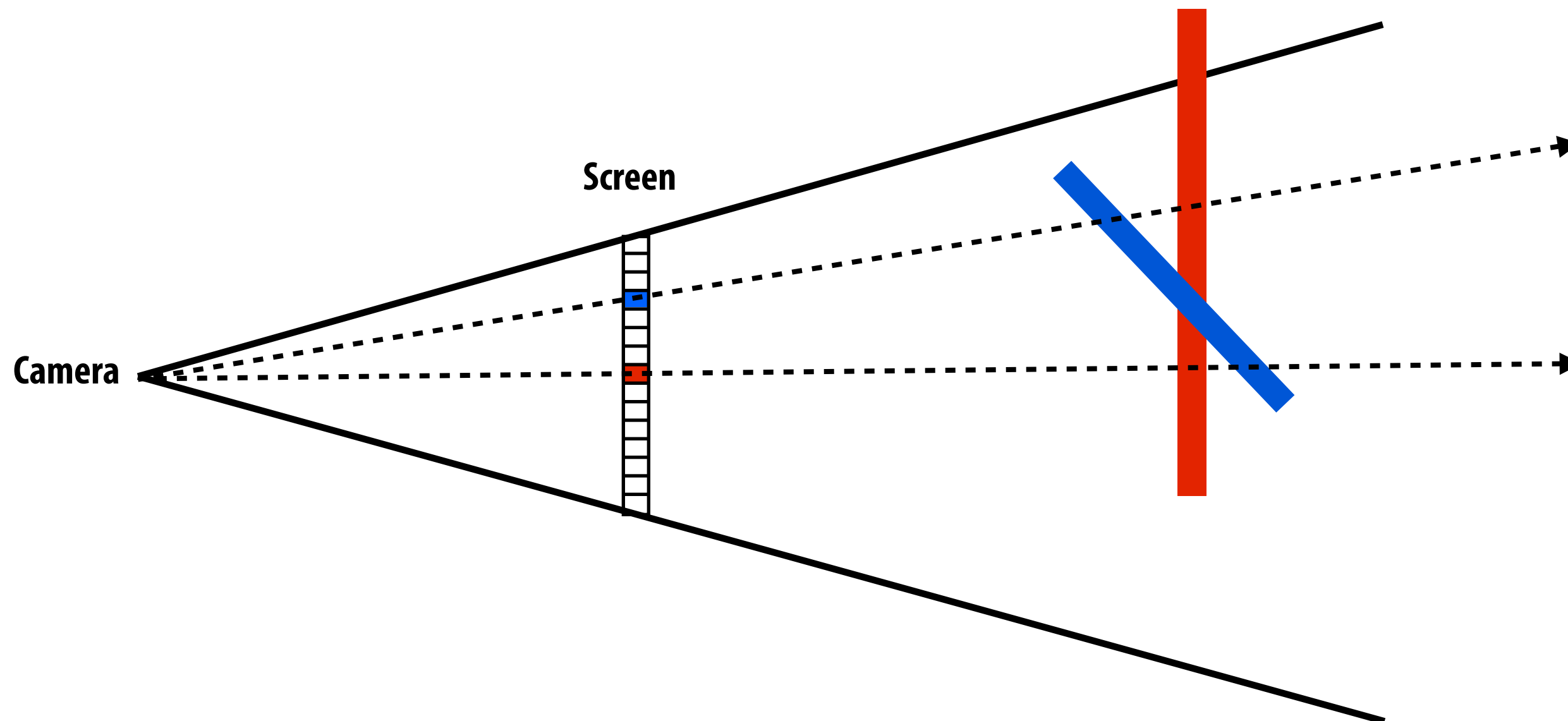
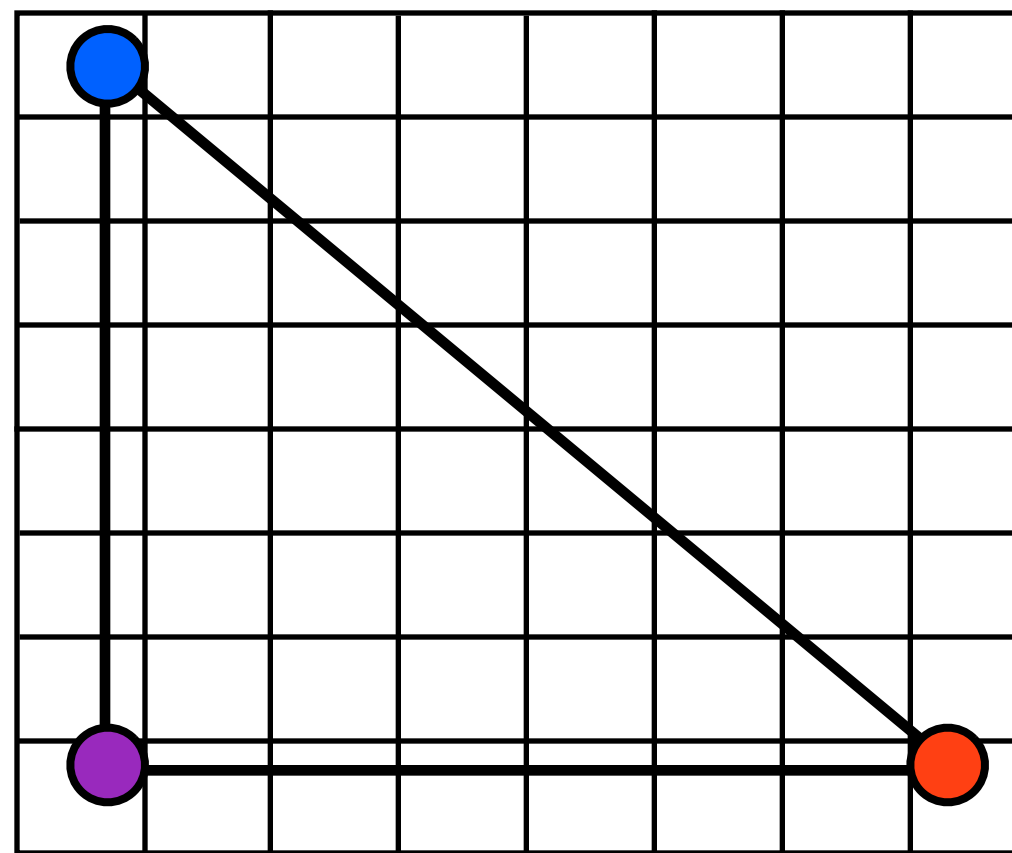


Image synthesis using the graphics pipeline

(As taught in graphics 101)

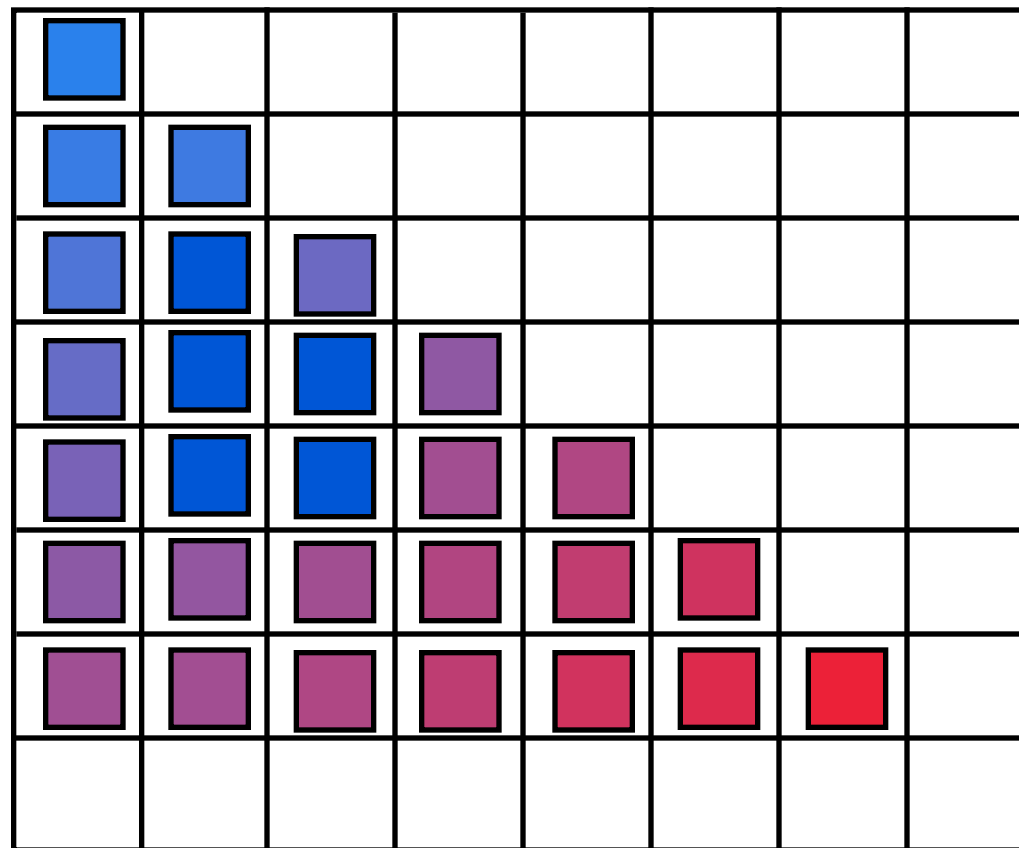
For each vertex, compute its projected position
(and other surface attributes, like normal, color)...



Then given a projected triangle...

Image synthesis by the graphics pipeline

(How it is taught in graphics 101)



```
struct my_fragment
{
    // application-defined attributes (opaque to pipeline)

    float3 normal;           // surface normal
    float2 texcoord1;        // texture coordinate for texture 1
    float2 texcoord2;        // texture coordinate for texture 2
    float3 worldPosition;    // world-space position of surface

    // pipeline-interpreted fields:

    int x, y;                // screen pixel position of fragment
    float depth;             // triangle depth for fragment
};
```

Rasterization converts the projected triangle into fragments.

Issue 1: determining coverage (what fragments should be generated?)

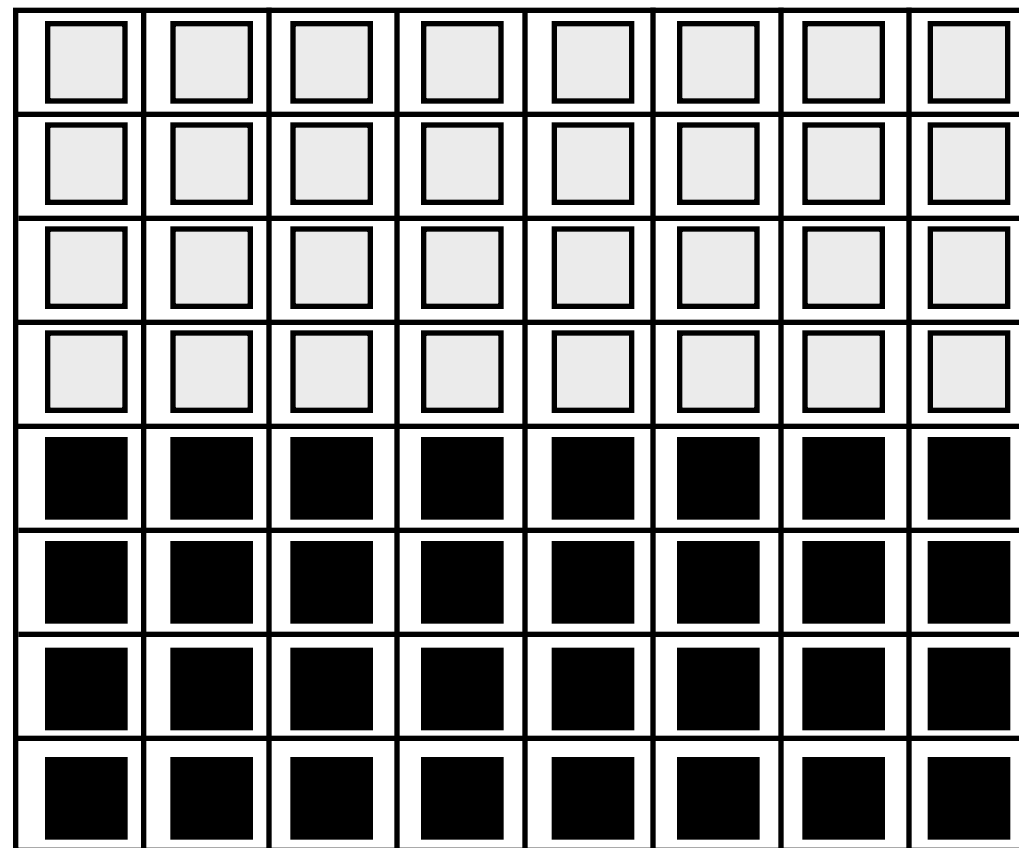
Issue 2: attribute interpolation (how to compute the value of surface attributes for each fragment?)

Image synthesis by the graphics pipeline

(How it is taught in graphics 101)

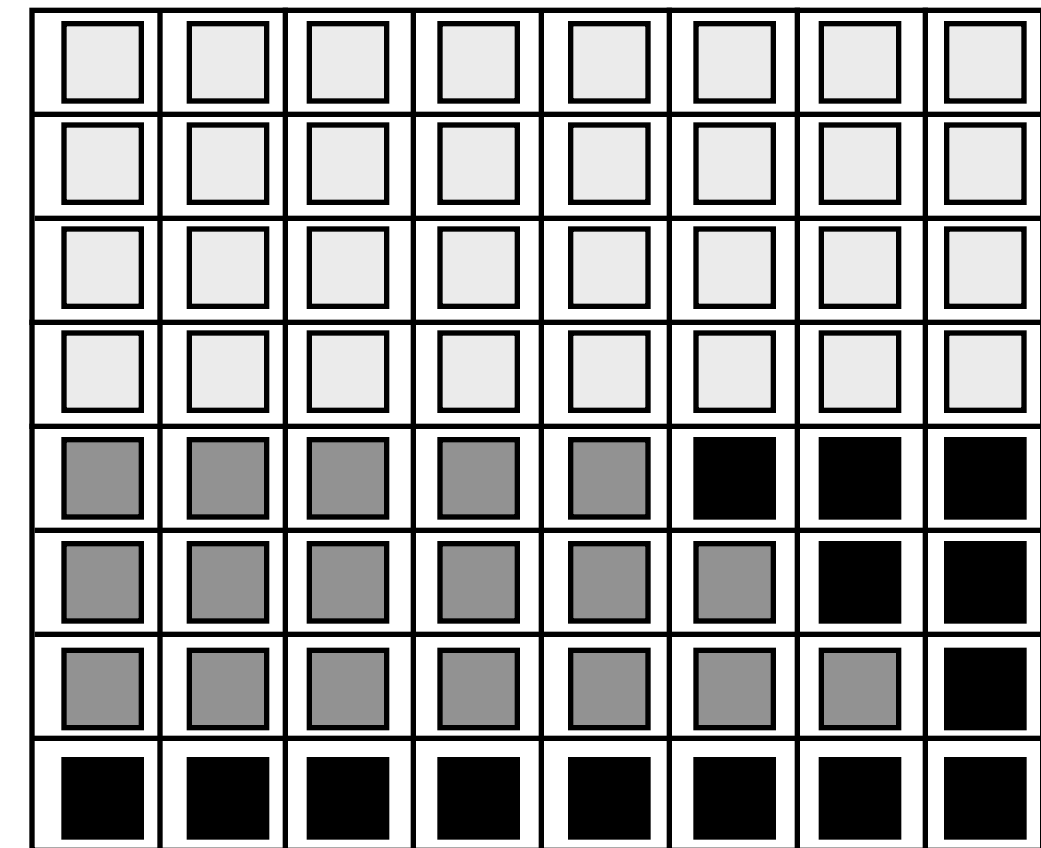
Z-buffer algorithm is used to determine the “closest” fragment at each pixel

Depth buffer before processing
fragments from new triangle



Near values = □ Far values = ■

Depth buffer after processing
fragments from new triangle



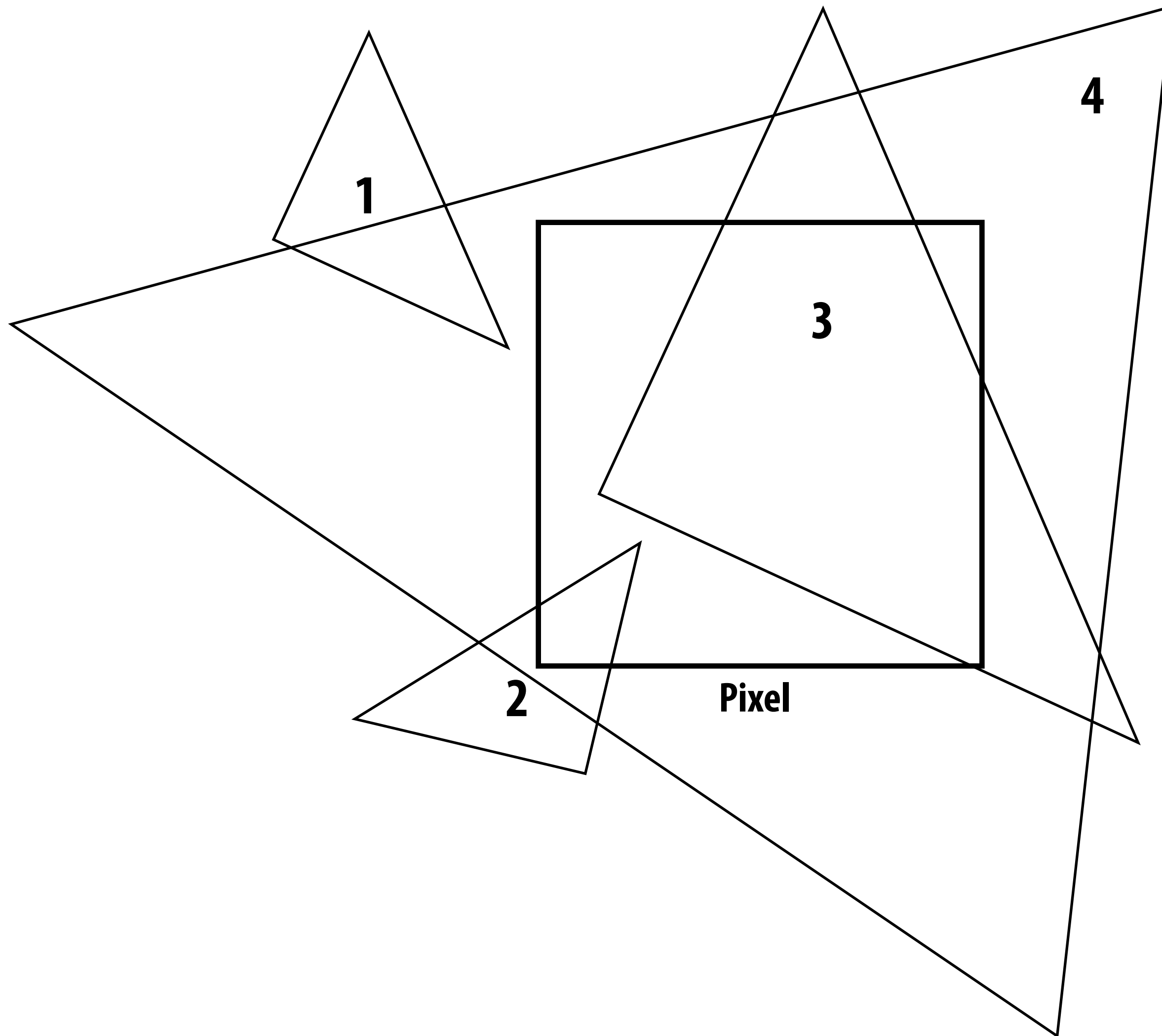
Let depth of new triangle = ■

```
bool pass_depth_test(float a, float b) {  
    return a < b;    // less-than predicate (the predicate is configurable in a modern pipeline)  
}
```

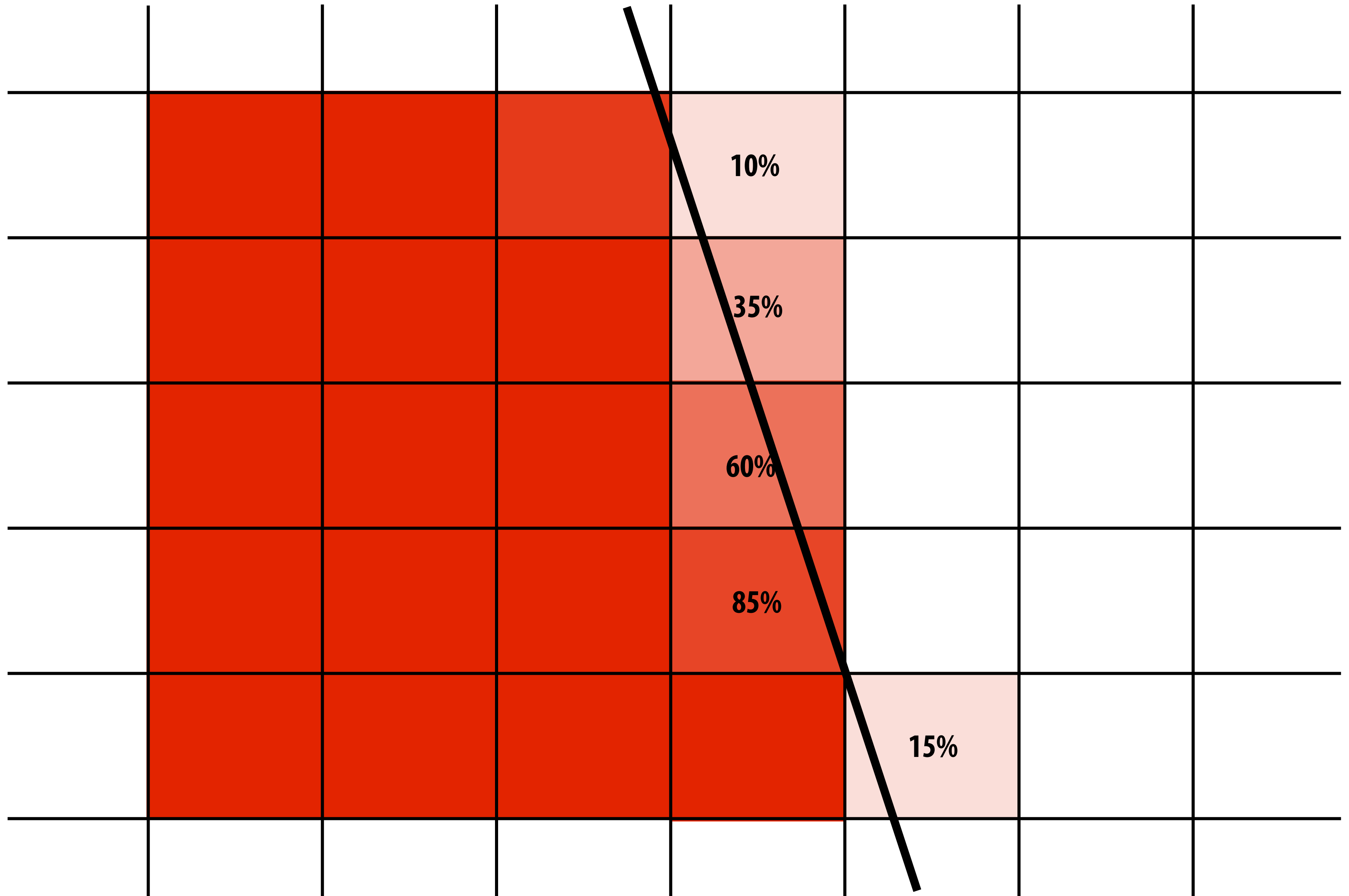
```
// “depth test”  
if (pass_depth_test(input_frag.depth, zbuffer[input_frag.x][input_frag.y])  
{  
    zbuffer[input_frag.x][input_frag.y] = input_frag.depth;  
    *** UPDATE COLOR BUFFER HERE AS WELL (NOT SHOWN) ***  
}
```

What does it mean for a pixel to be covered by a triangle?

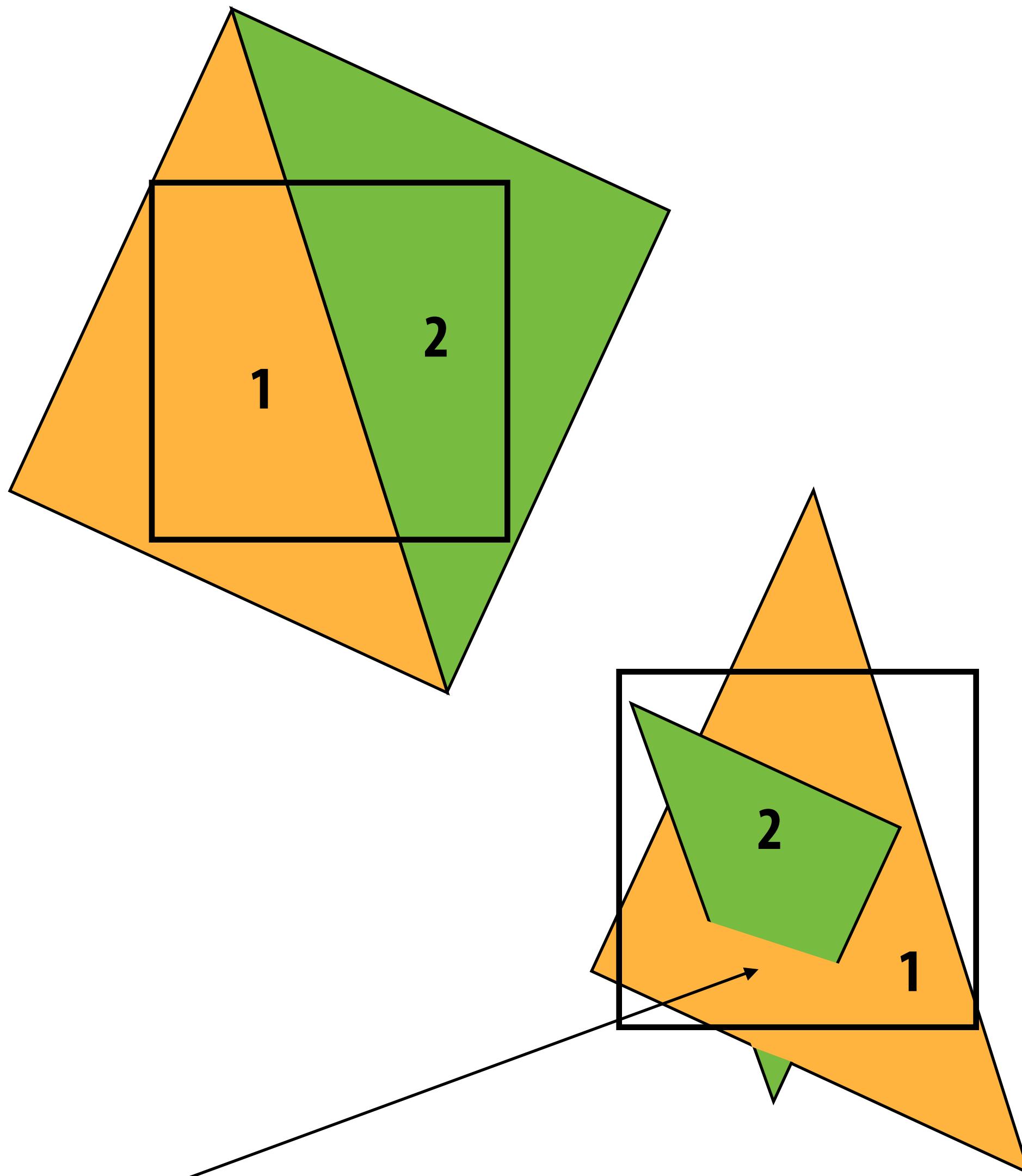
Question: which triangles “cover” this pixel?



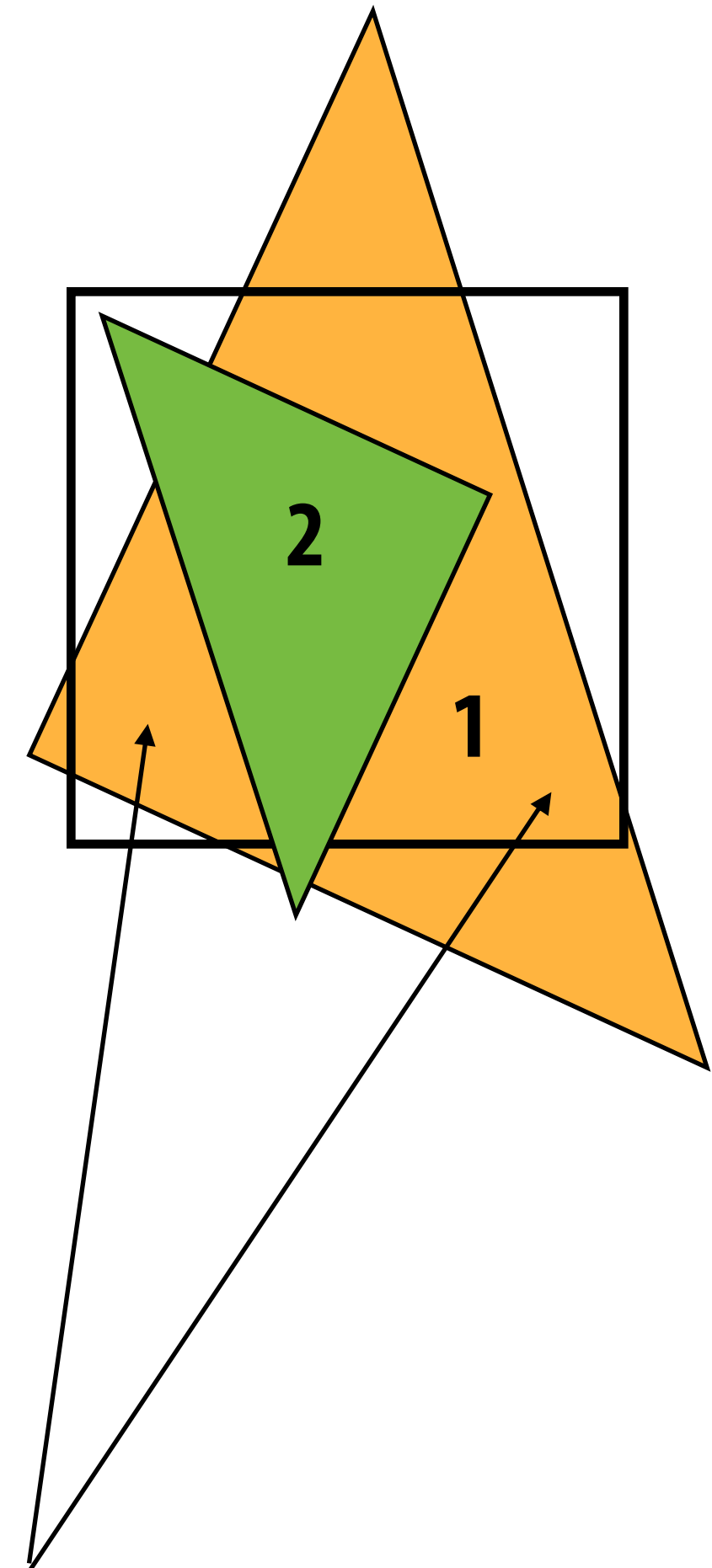
One option: analytically compute fraction of pixel covered by triangle



Analytical schemes get tricky when considering occlusion



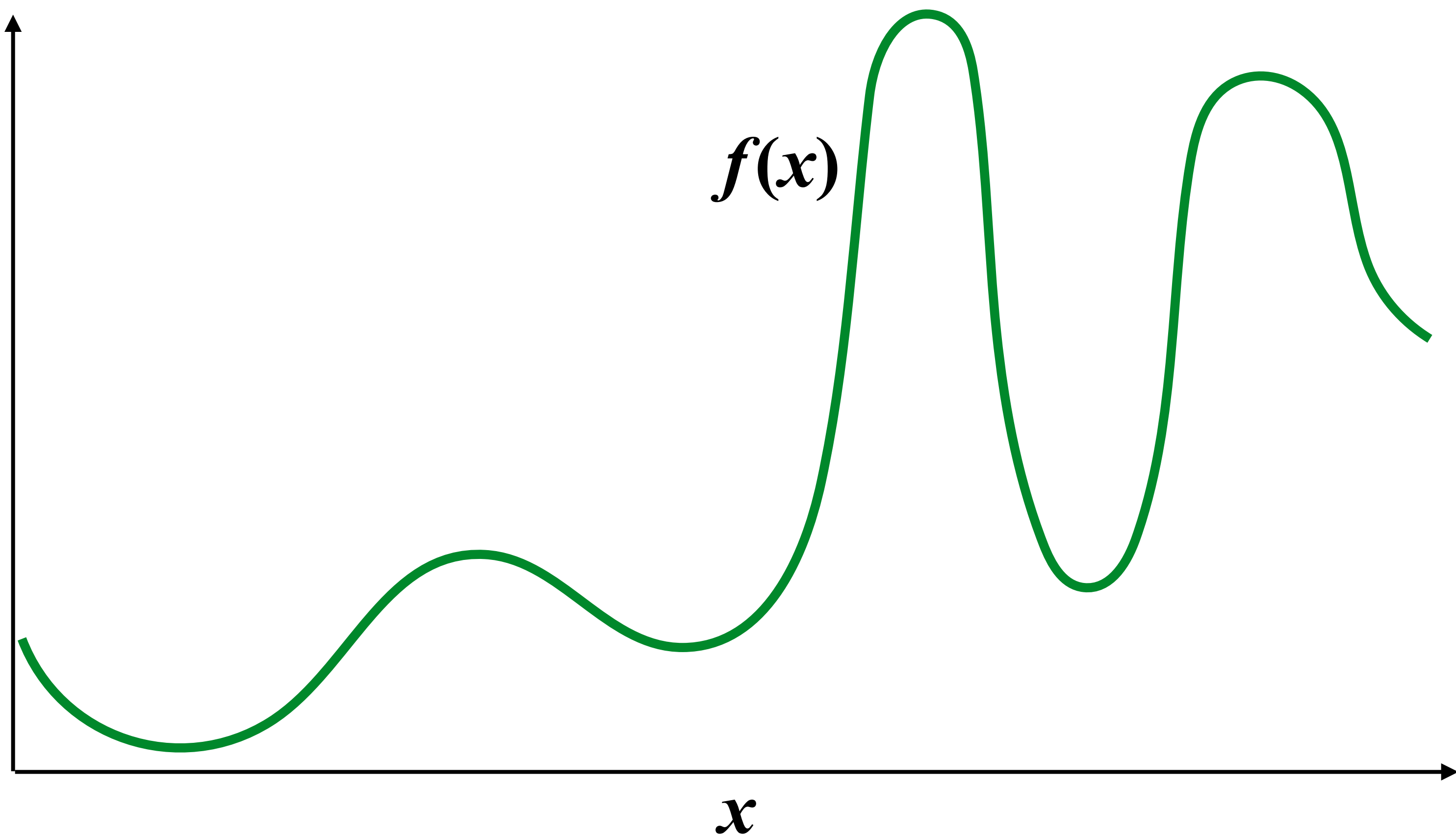
Interpenetration: even worse



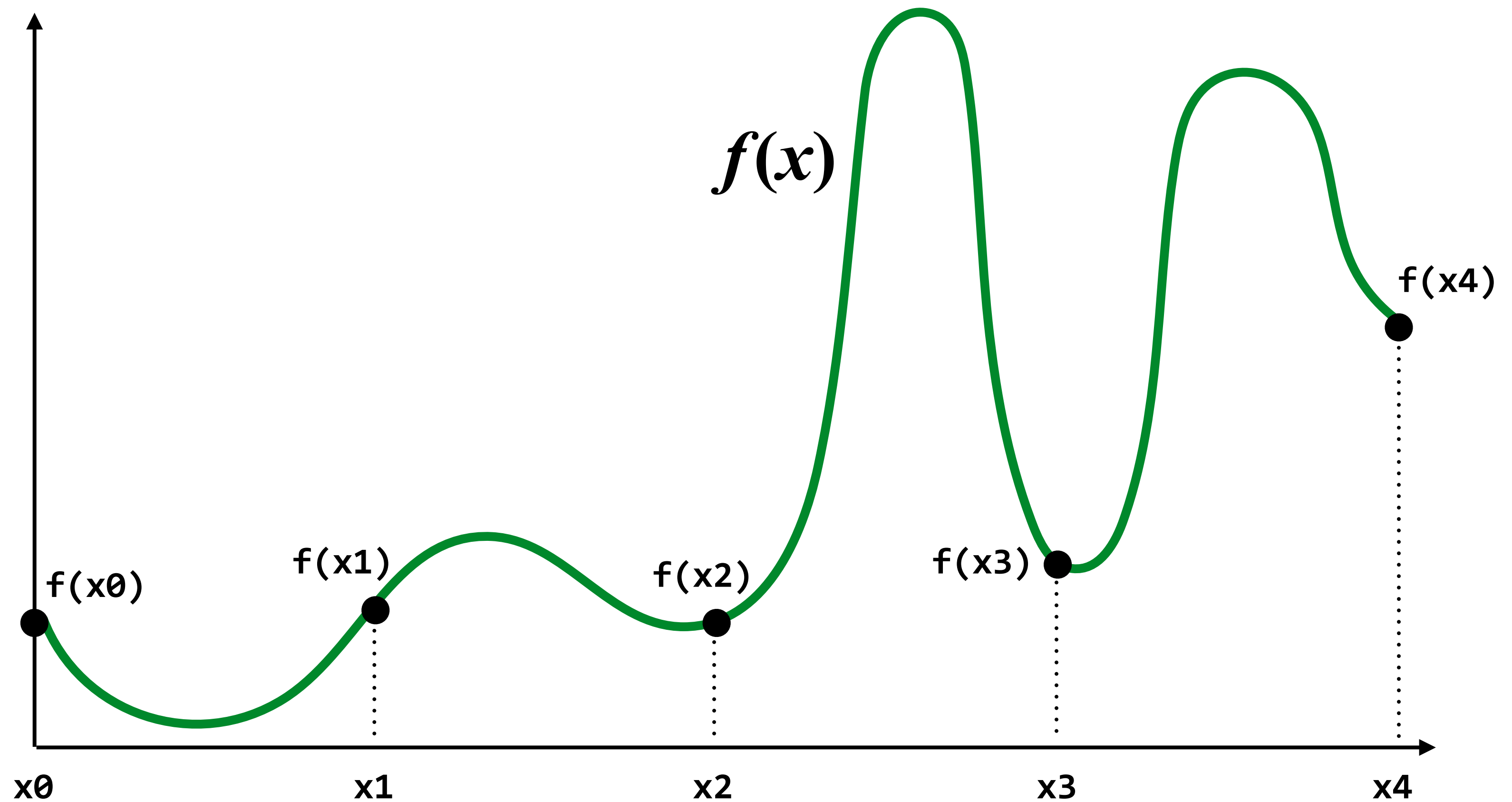
Two regions of [1] contribute to pixel. One of these regions is not even convex.

Sampling 101

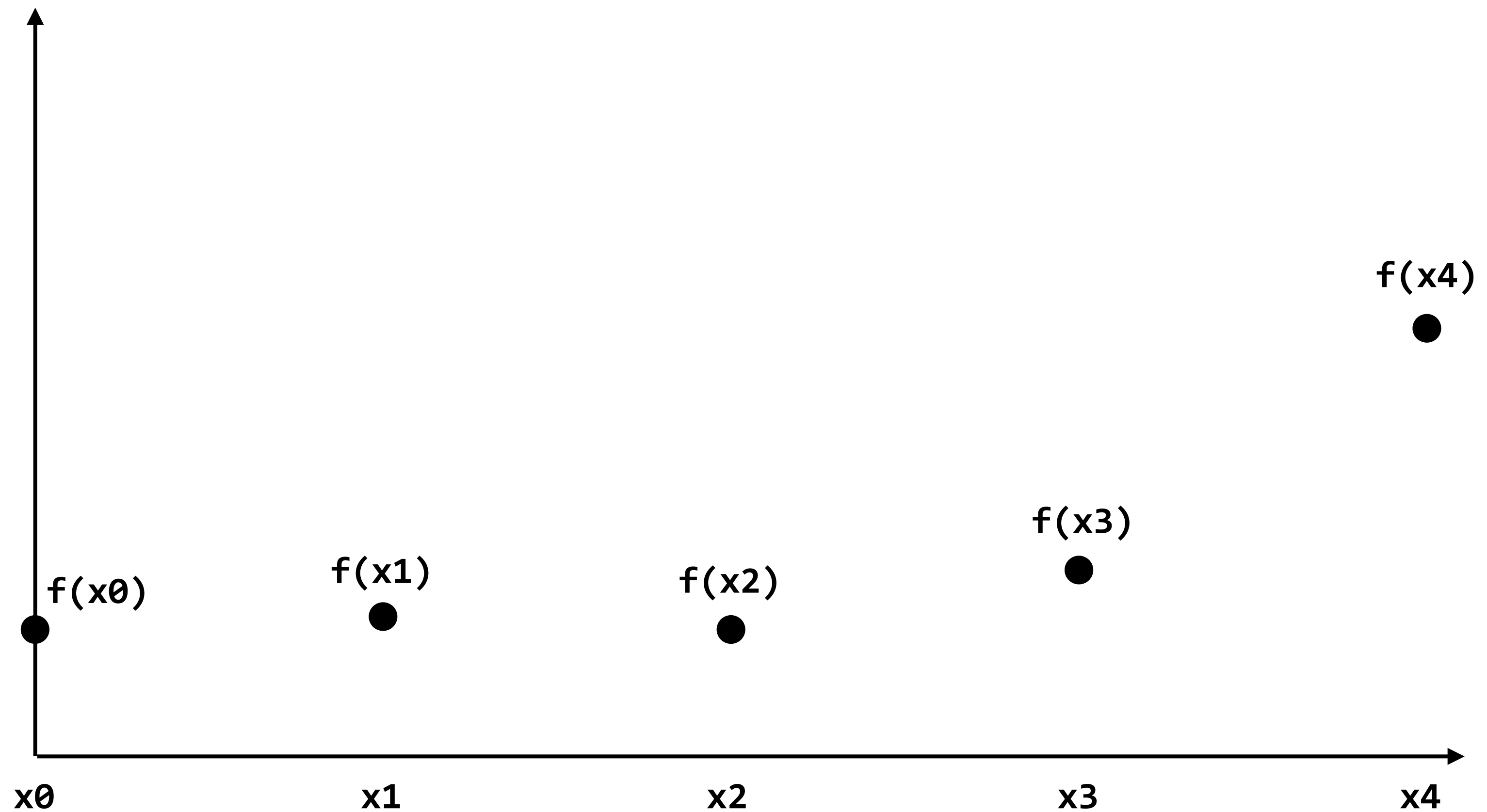
Continuous 1D signal



Measurements of 1D signal (5 measurements)

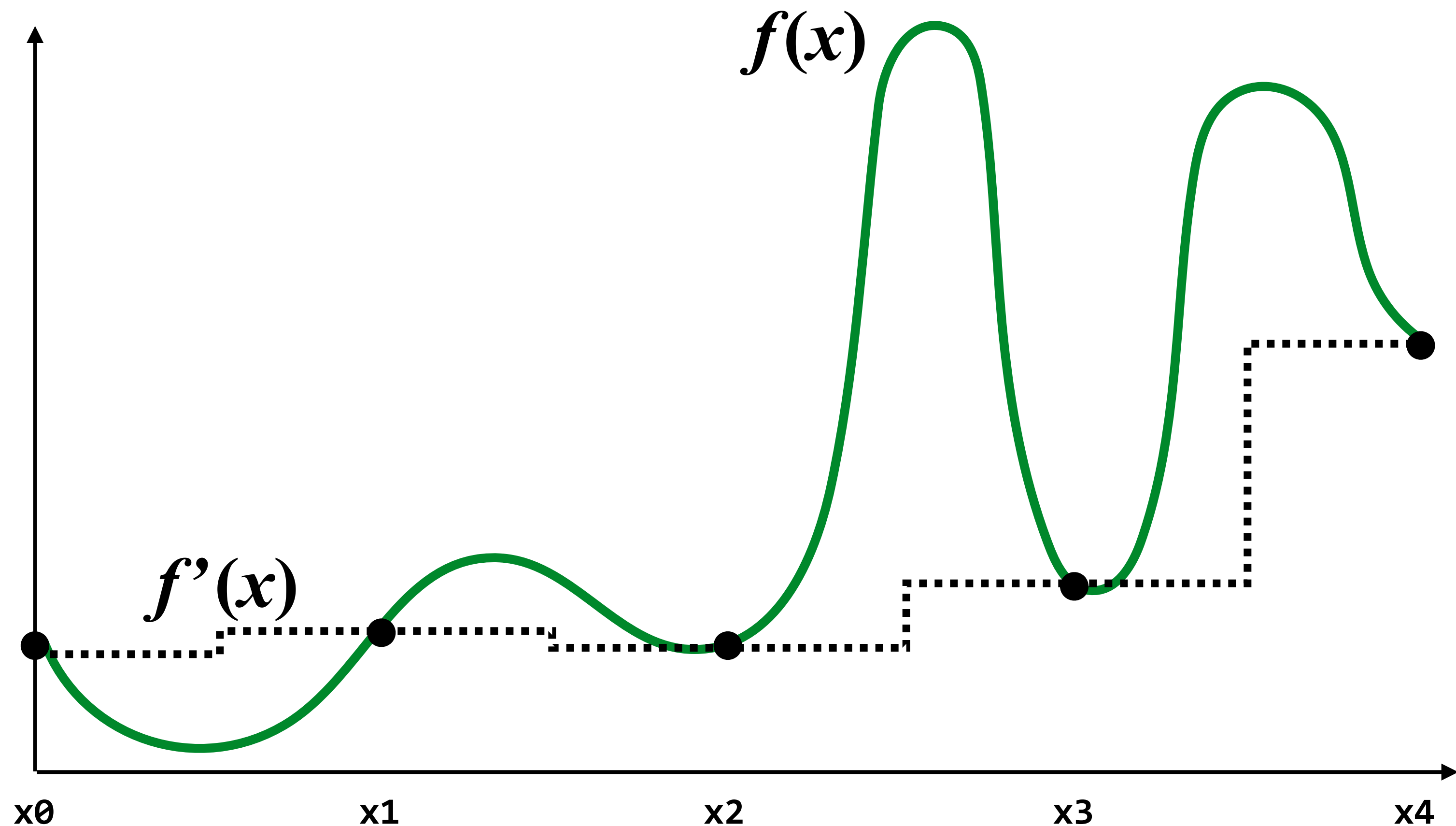


Reconstructing the signal

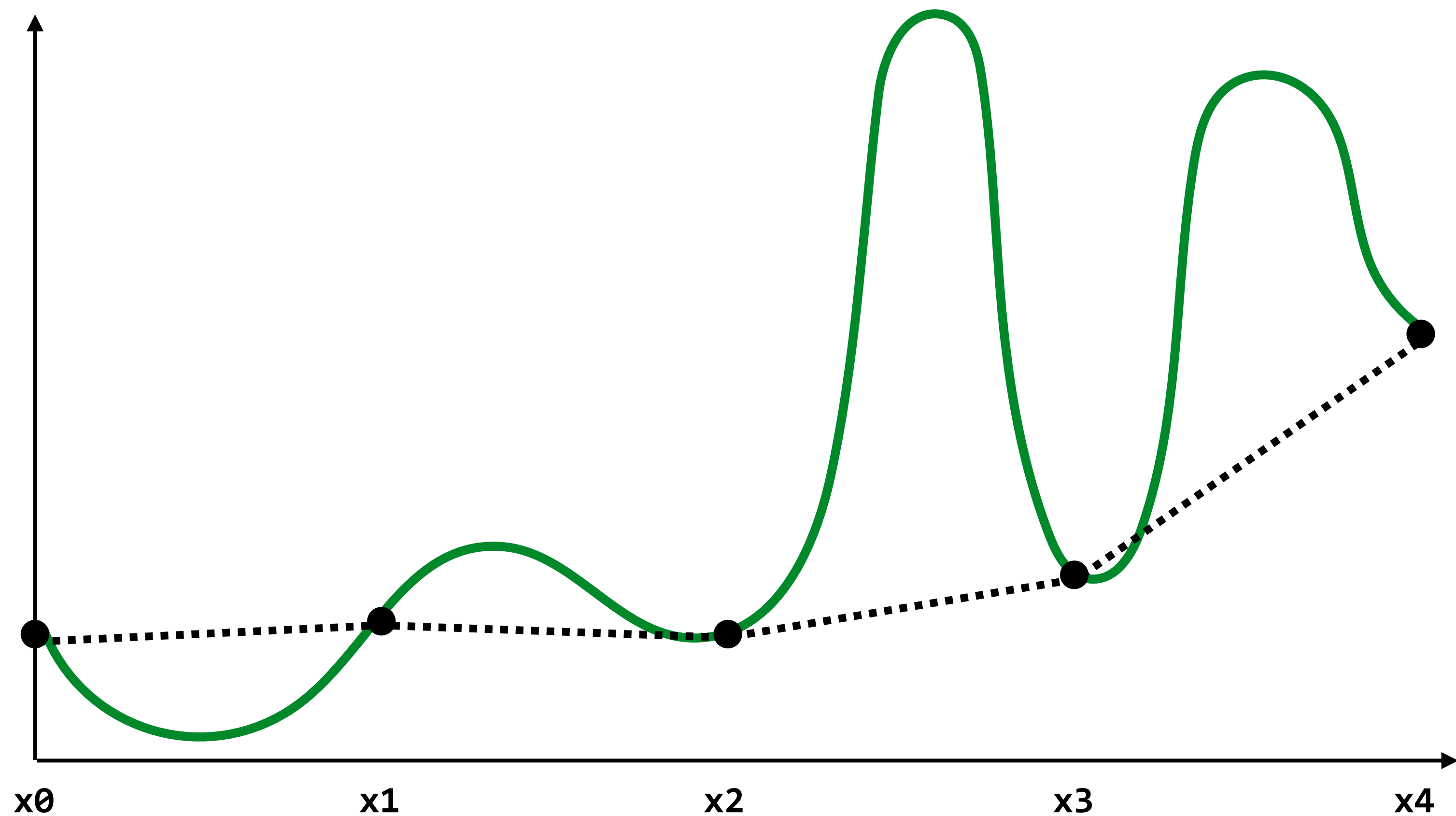


Piecewise constant approximation

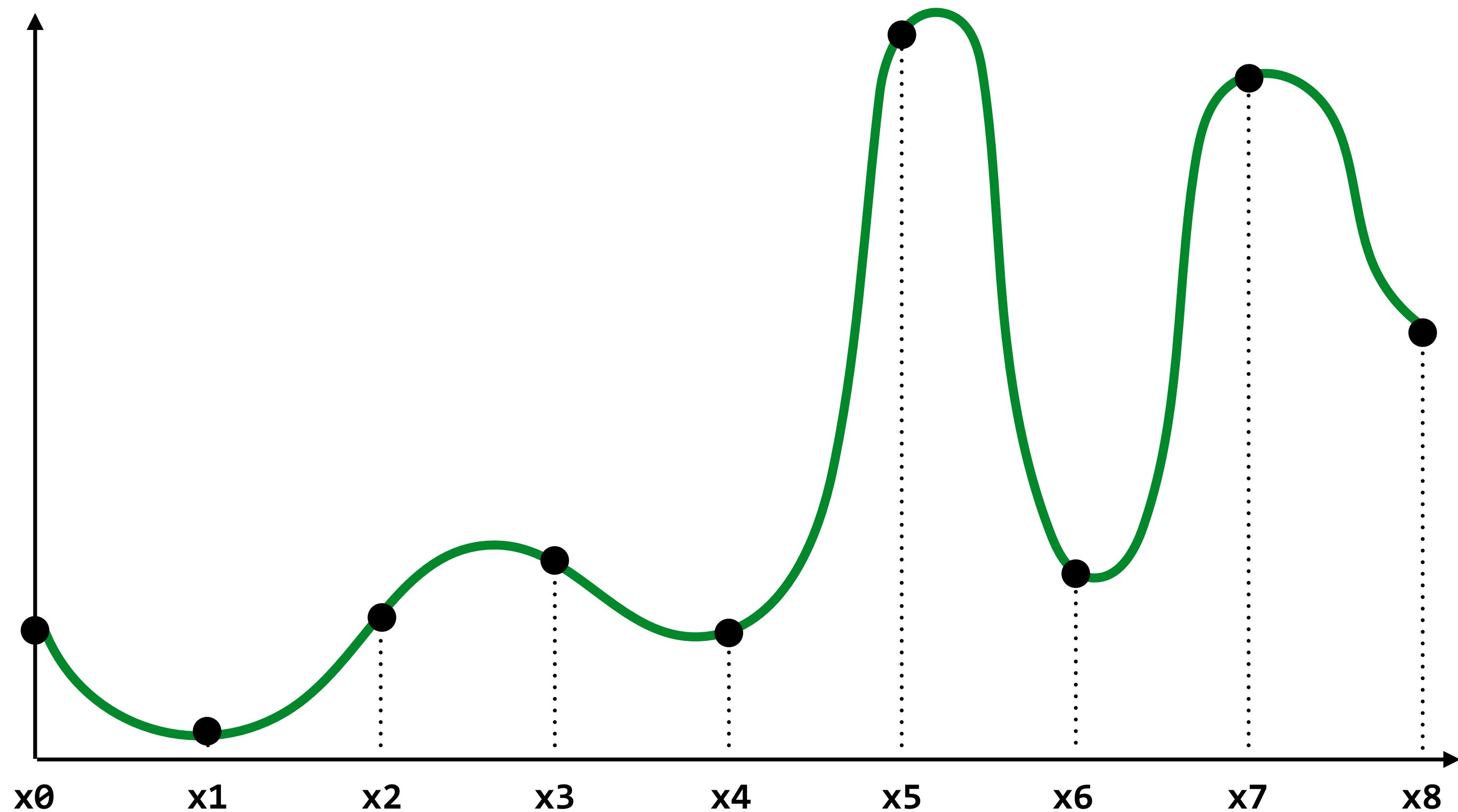
$f'(x)$ approximates $f(x)$



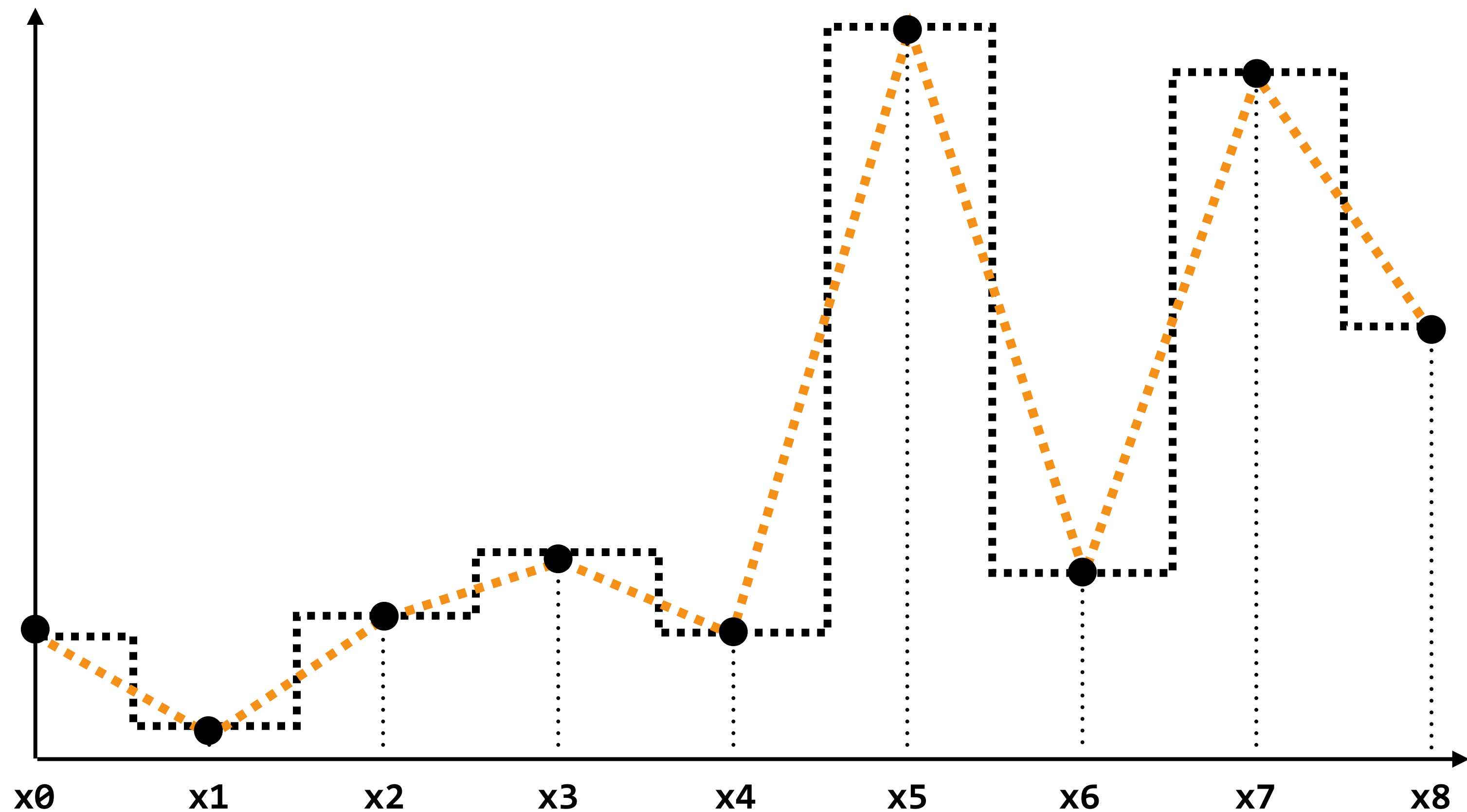
Piecewise linear approximation



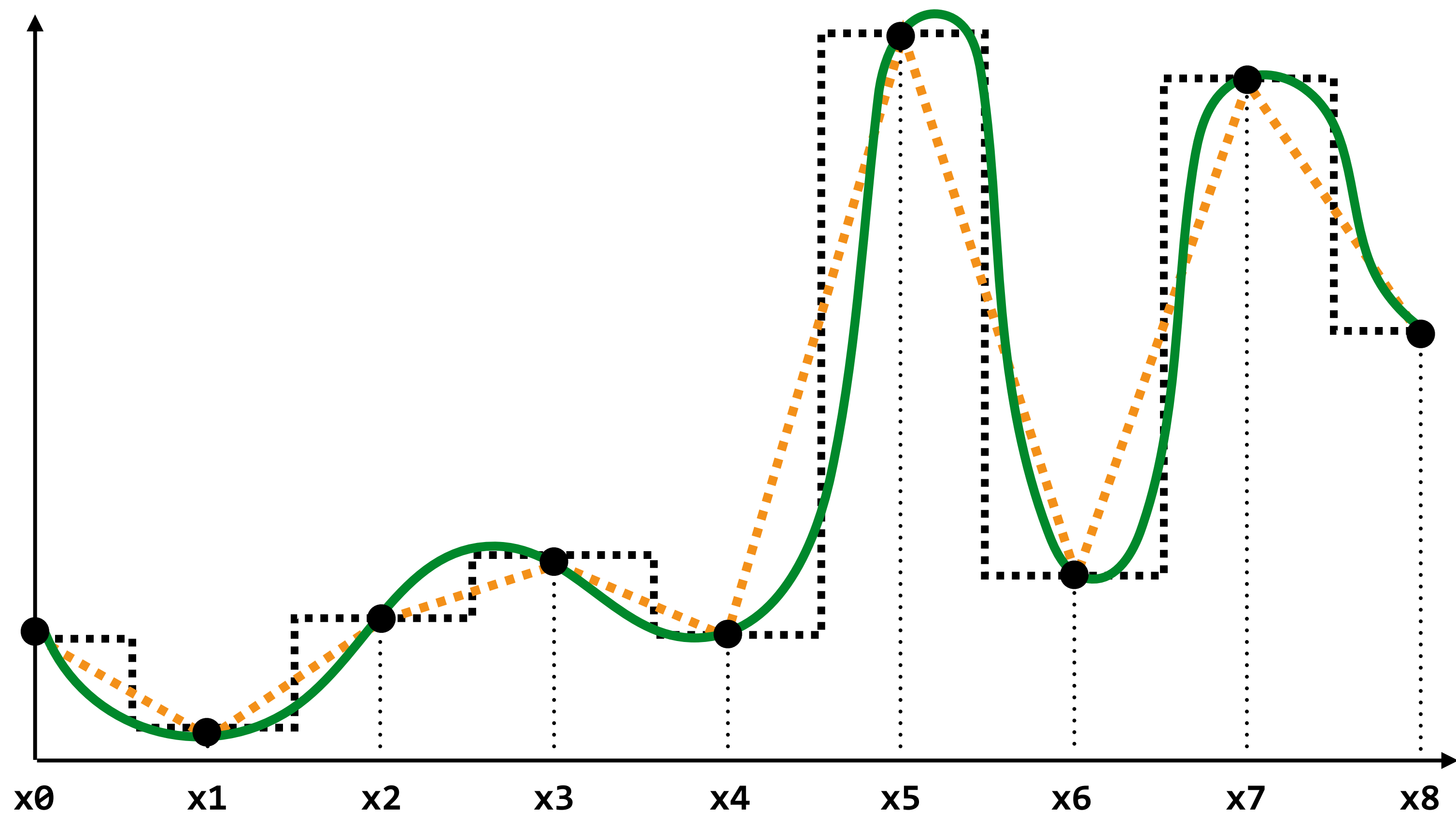
How can we represent the signal more accurately?



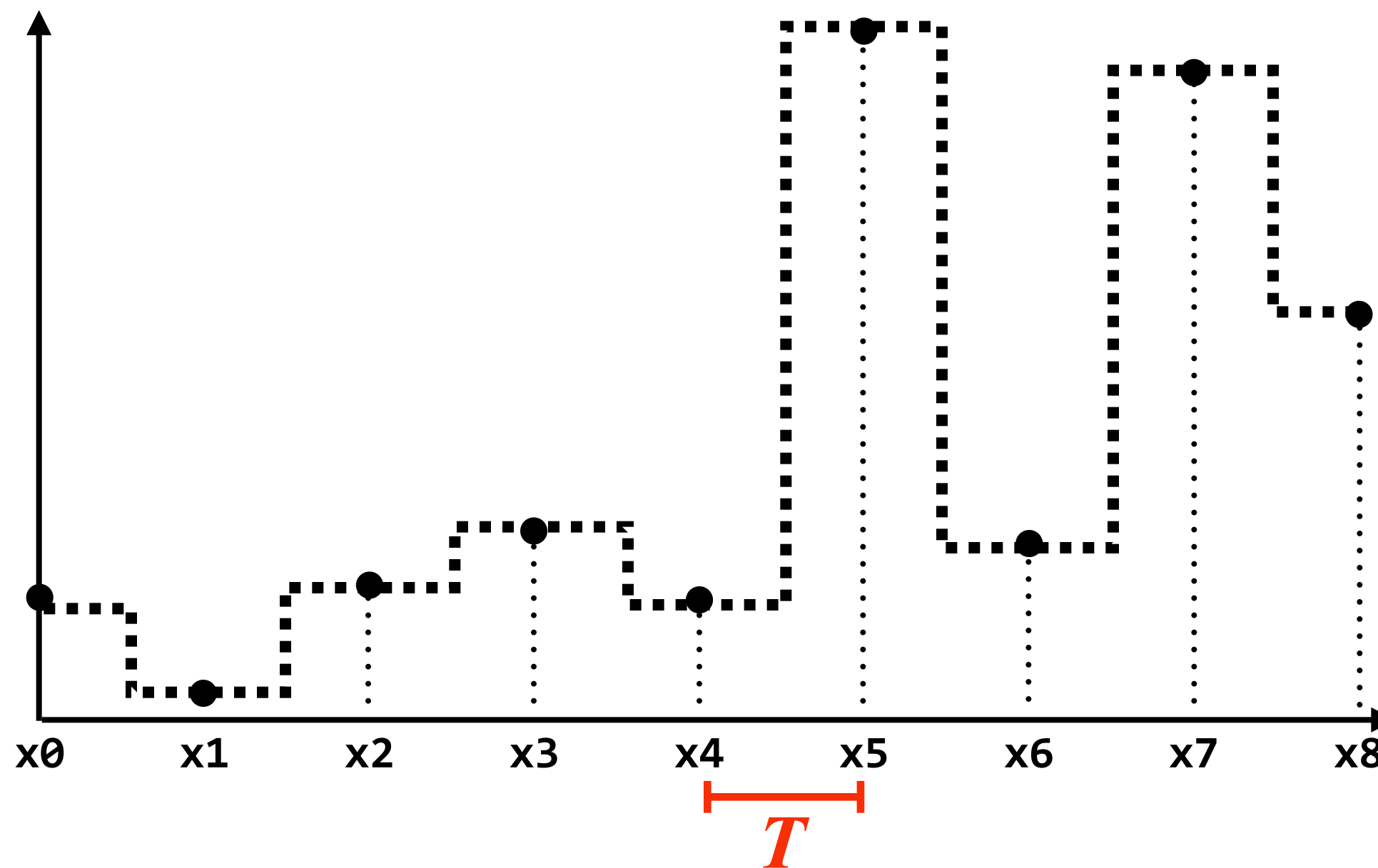
Reconstruction from denser sampling



Reconstruction from denser sampling



Reconstruction as convolution (box filter)



Sampled signal

$$g(x) = \sum_{i=0}^8 f(iT) \delta(x - iT)$$

Question: what if we define reconstruction filter to be:

Reconstruction filter (box)

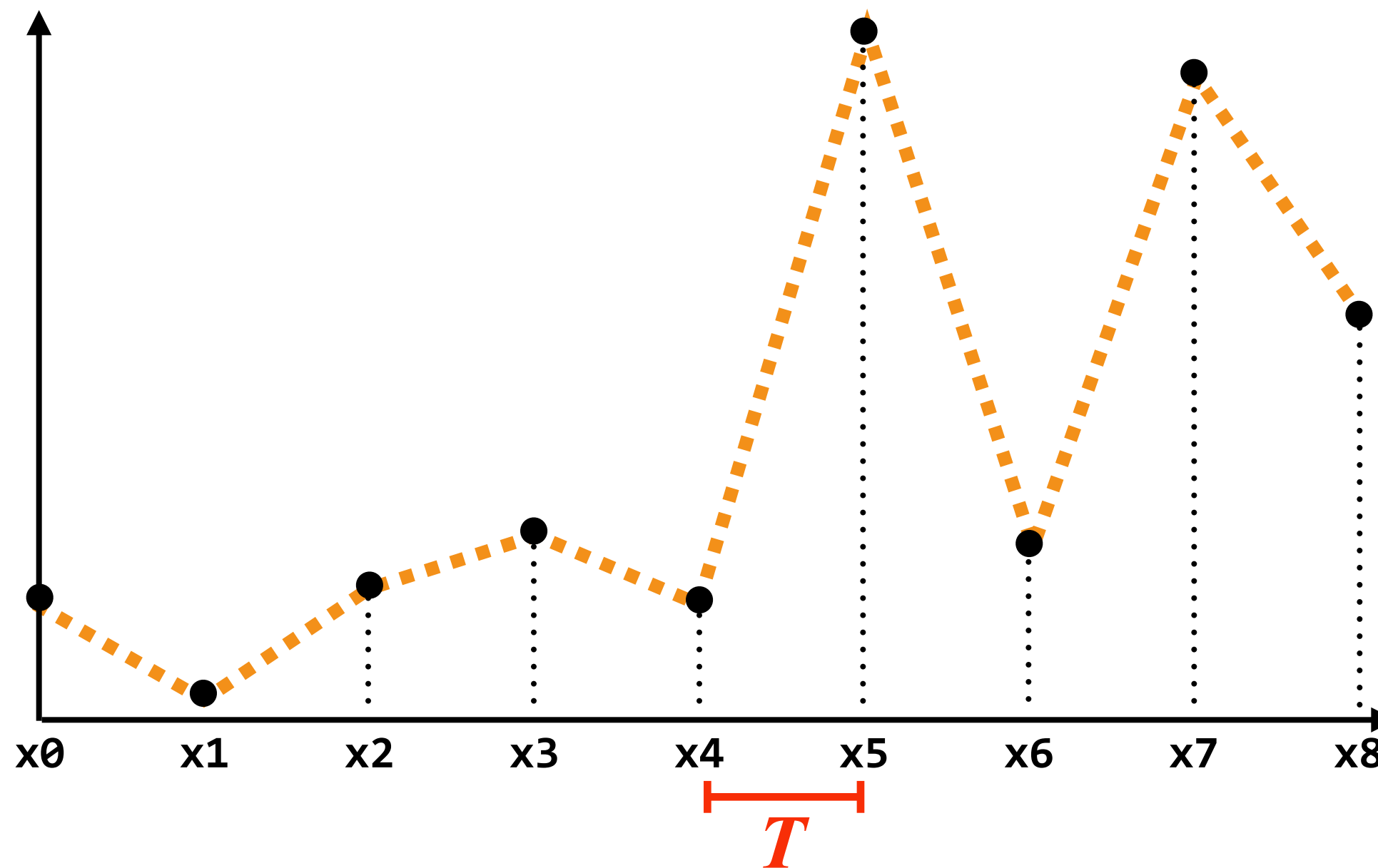
$$h(x) = \begin{cases} 1, & \text{if } |x| \leq T/2 \\ 0, & \text{if } |x| > T/2 \end{cases}$$

$$h(x) = \begin{cases} 1/3 & \text{if } |x| \leq T \\ 0 & \text{if } |x| > T \end{cases}$$

**Constructed signal
(Chooses nearest sample)**

$$f'(x) = \int_{-\infty}^{\infty} h(x') g(x - x') dx' = \int_{-T/2}^{T/2} \sum_{i=0}^8 f(iT) \delta(x - x' - iT) dx'$$

Reconstruction as convolution (triangle filter)



Sampled signal

$$g(x) = \sum_{i=0}^8 f(iT) \delta(x - iT)$$

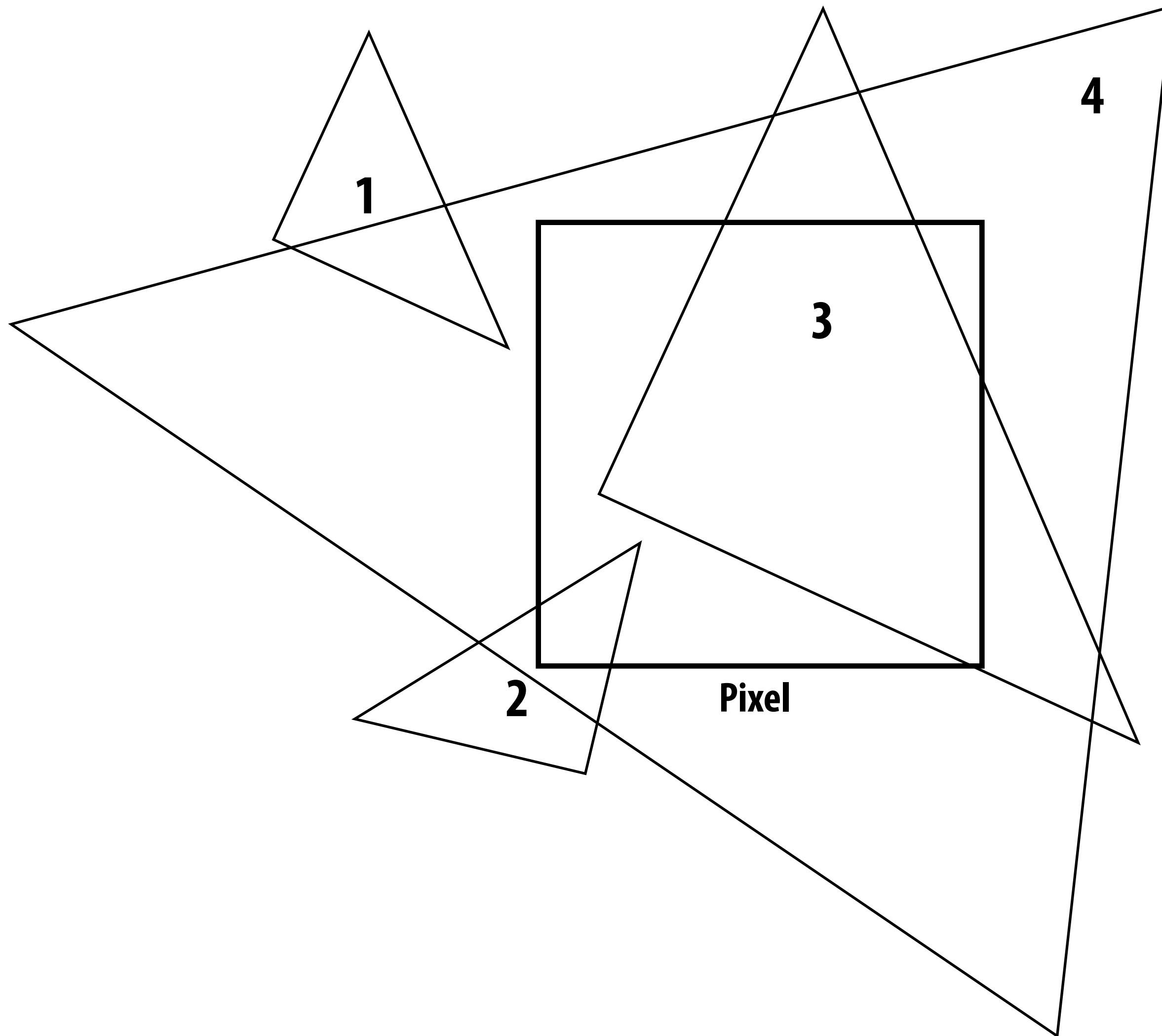
**Triangle reconstruction filter:
(yields linear interpolation of
samples)**

$$h(x) = \begin{cases} 1 - \frac{|x|}{T} & \text{if } |x| \leq T \\ 0 & \text{if } |x| > T \end{cases}$$

Coverage

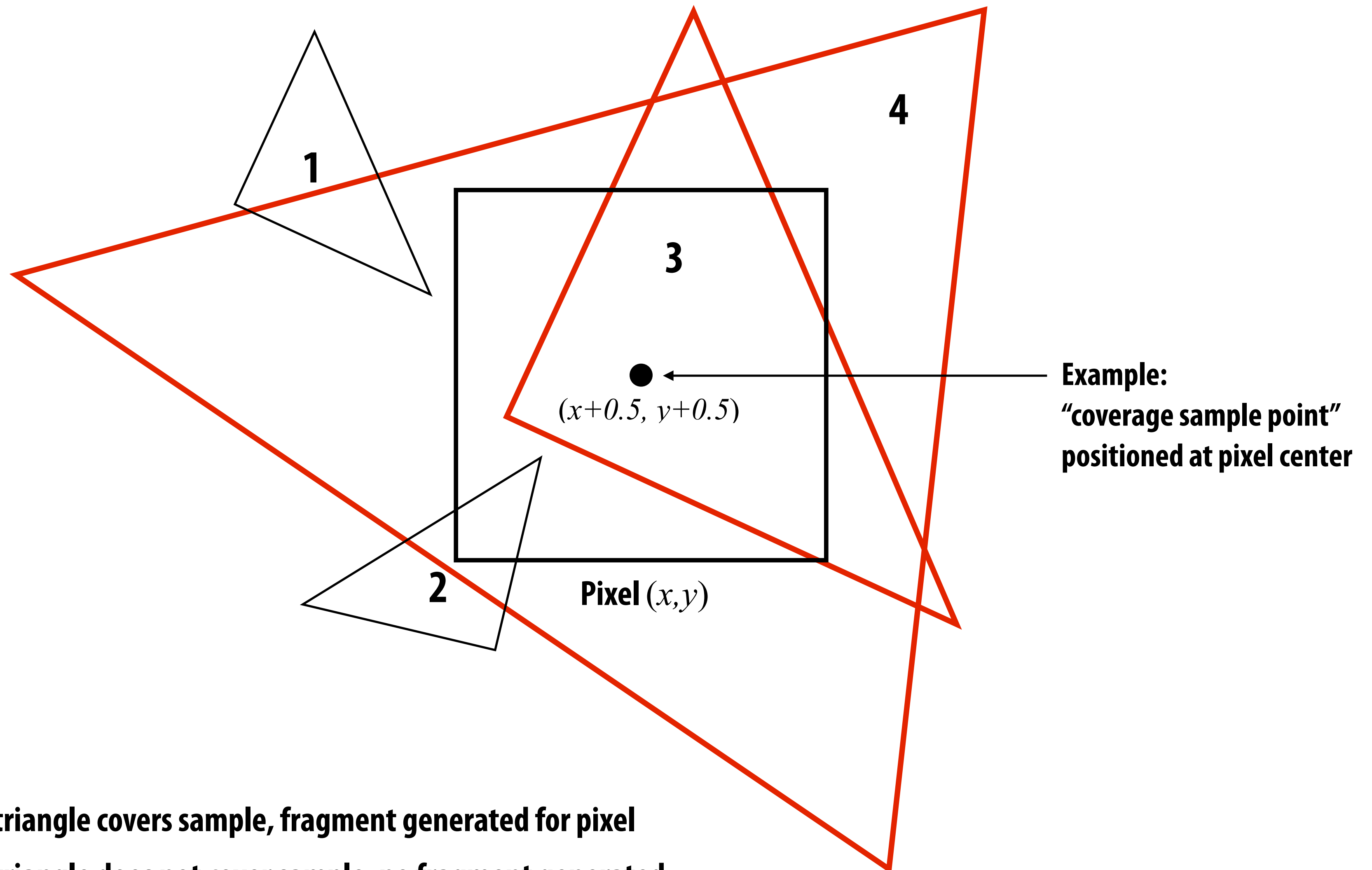
What does it mean for a pixel to be covered by a triangle?

Question: which triangles “cover” this pixel?



The rasterizer estimates triangle-screen coverage using point sampling

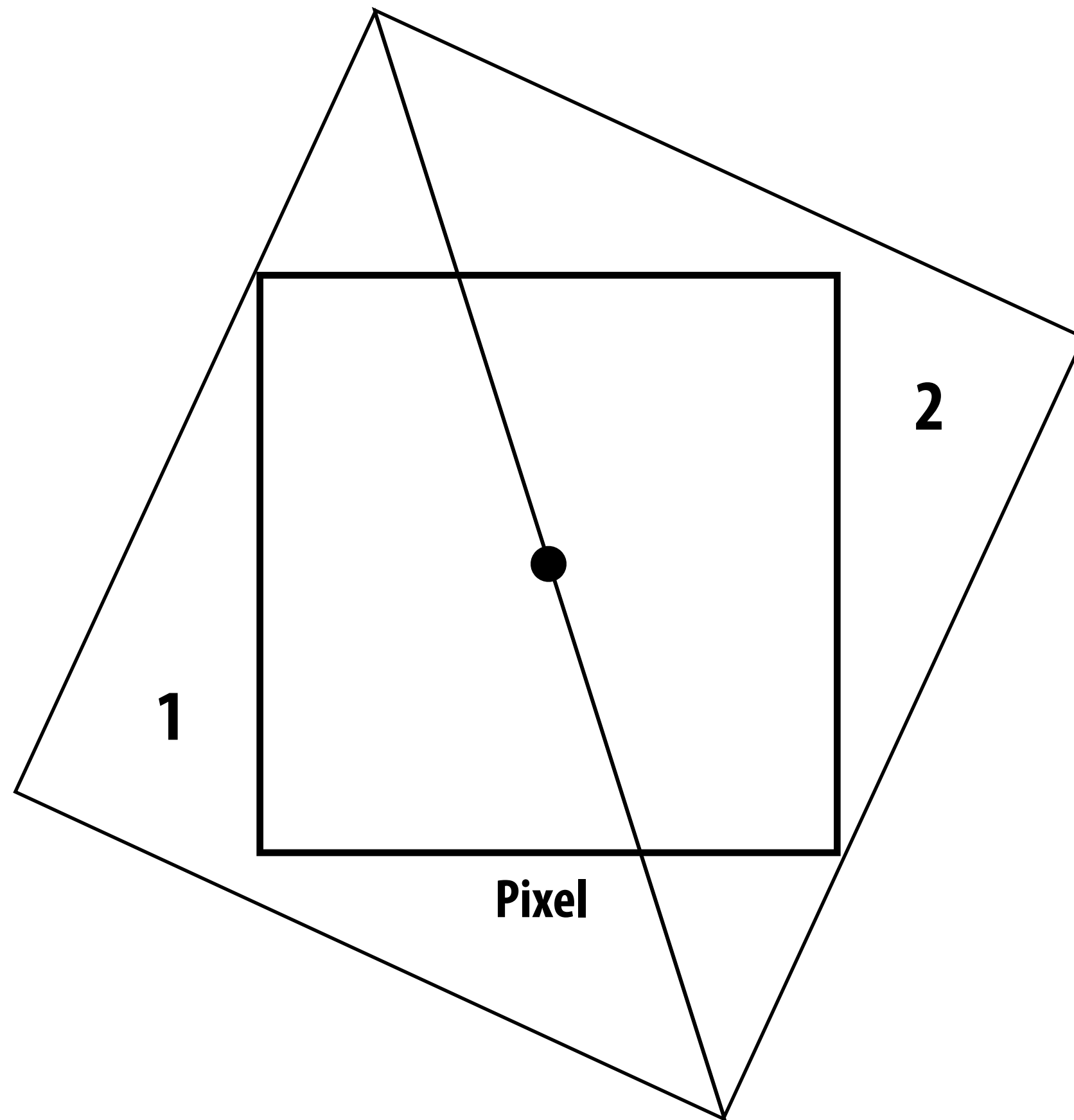
It samples the binary function: $\text{coverage}(x,y)$



-  = triangle covers sample, fragment generated for pixel
-  = triangle does not cover sample, no fragment generated

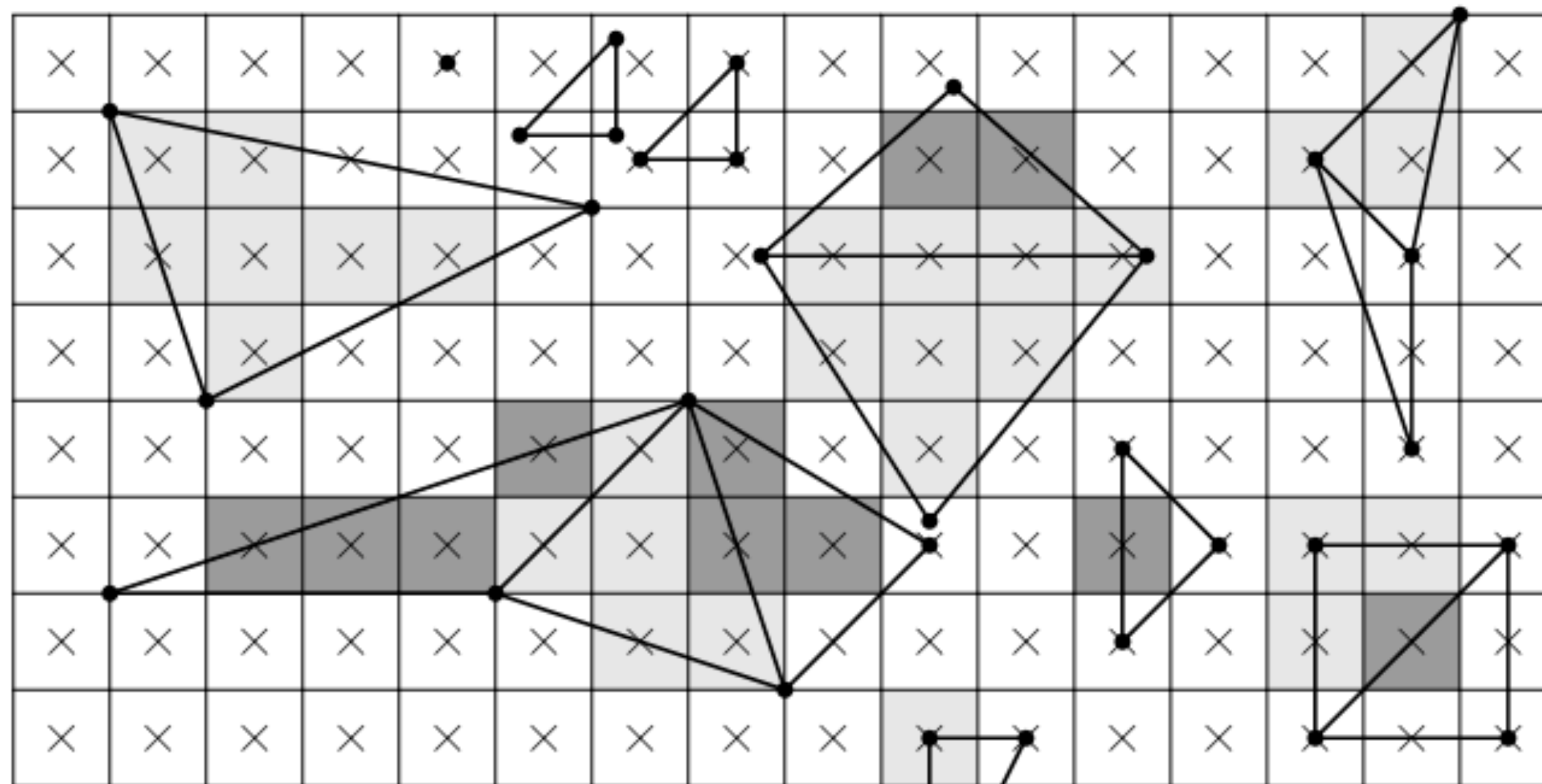
Edge cases (literally)

Is fragment generated for triangle 1? for triangle 2?



Edge rules

- **Direct3D rules:** when edge falls directly on sample, sample classified as within triangle if the edge is a “top edge” or “left edge”
 - **Top edge:** horizontal edge that is above all other edges
 - **Left edge:** an edge that is not exactly horizontal and is on the left side of the triangle. (triangle can have one or two left edges)



Source: Direct3D Programming Guide, Microsoft



Pixel
(cross = center; x,y @ 0.5)

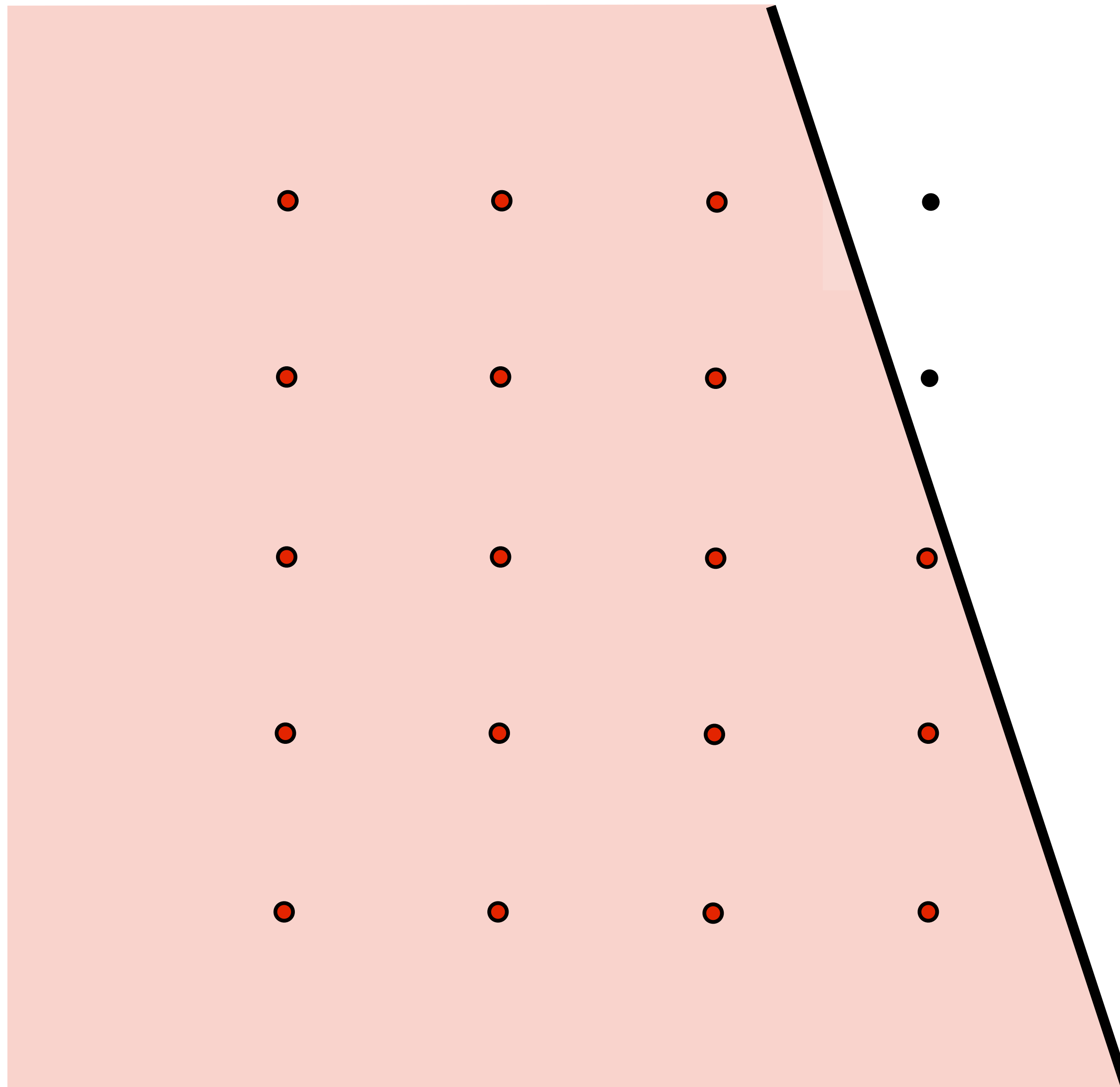


Triangle

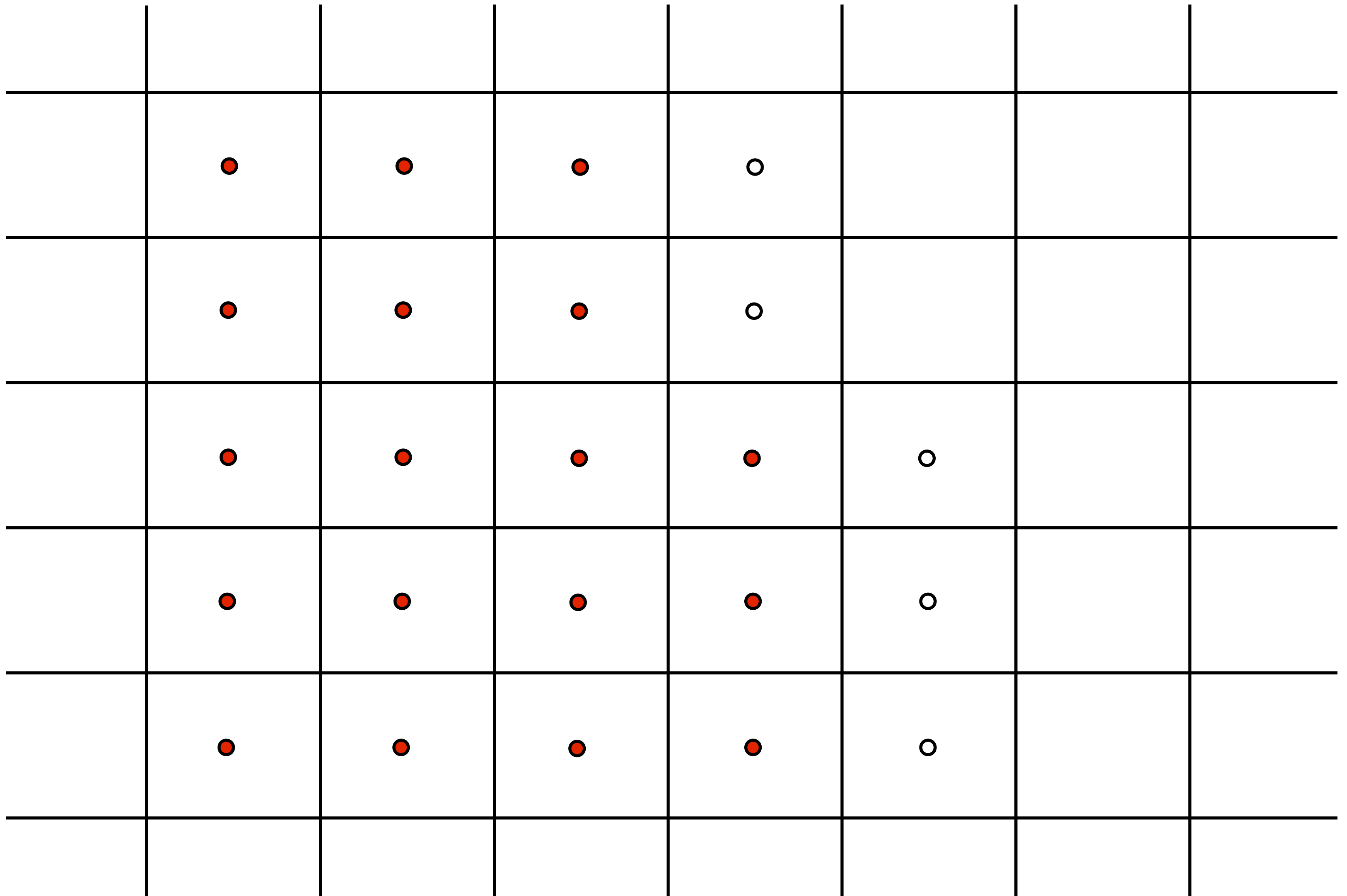


Covered
Pixels

Results of sampling triangle coverage

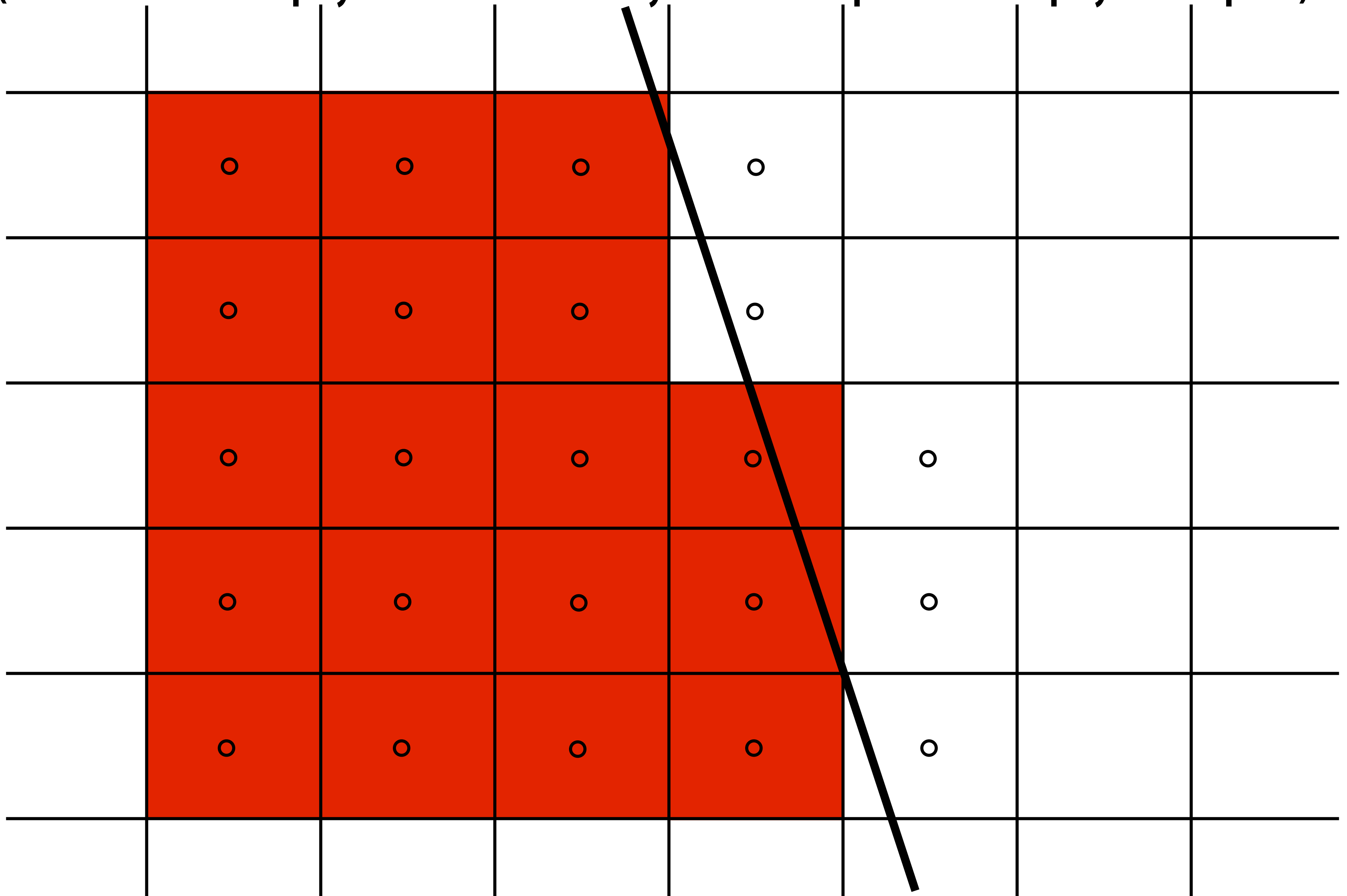


Results of sampling (red dots = covered)



Reconstruction with box filter (aliased edges)

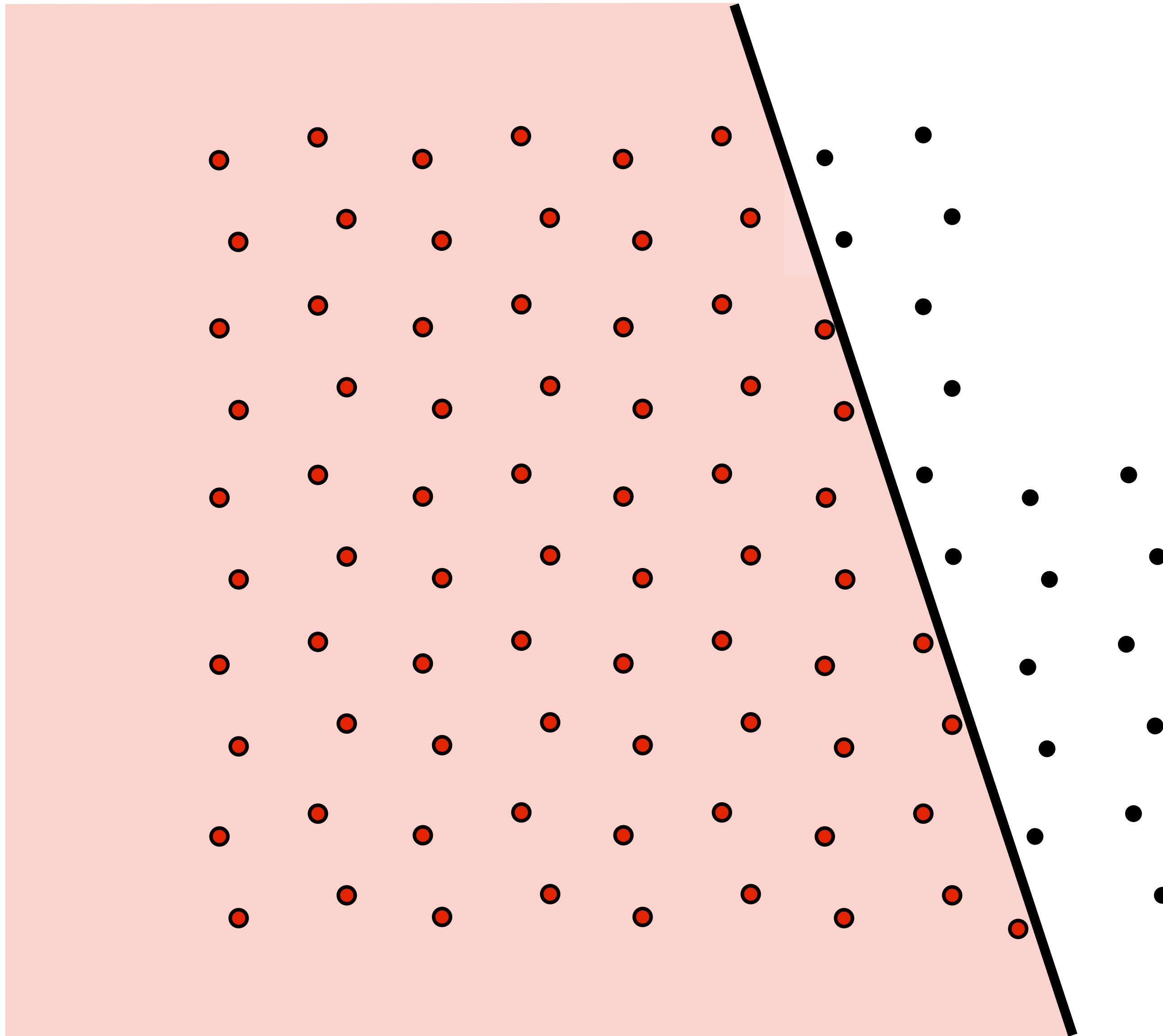
(consider this the displayed result: not actually true since a “pixel” on a display is not square)



Problem: aliasing

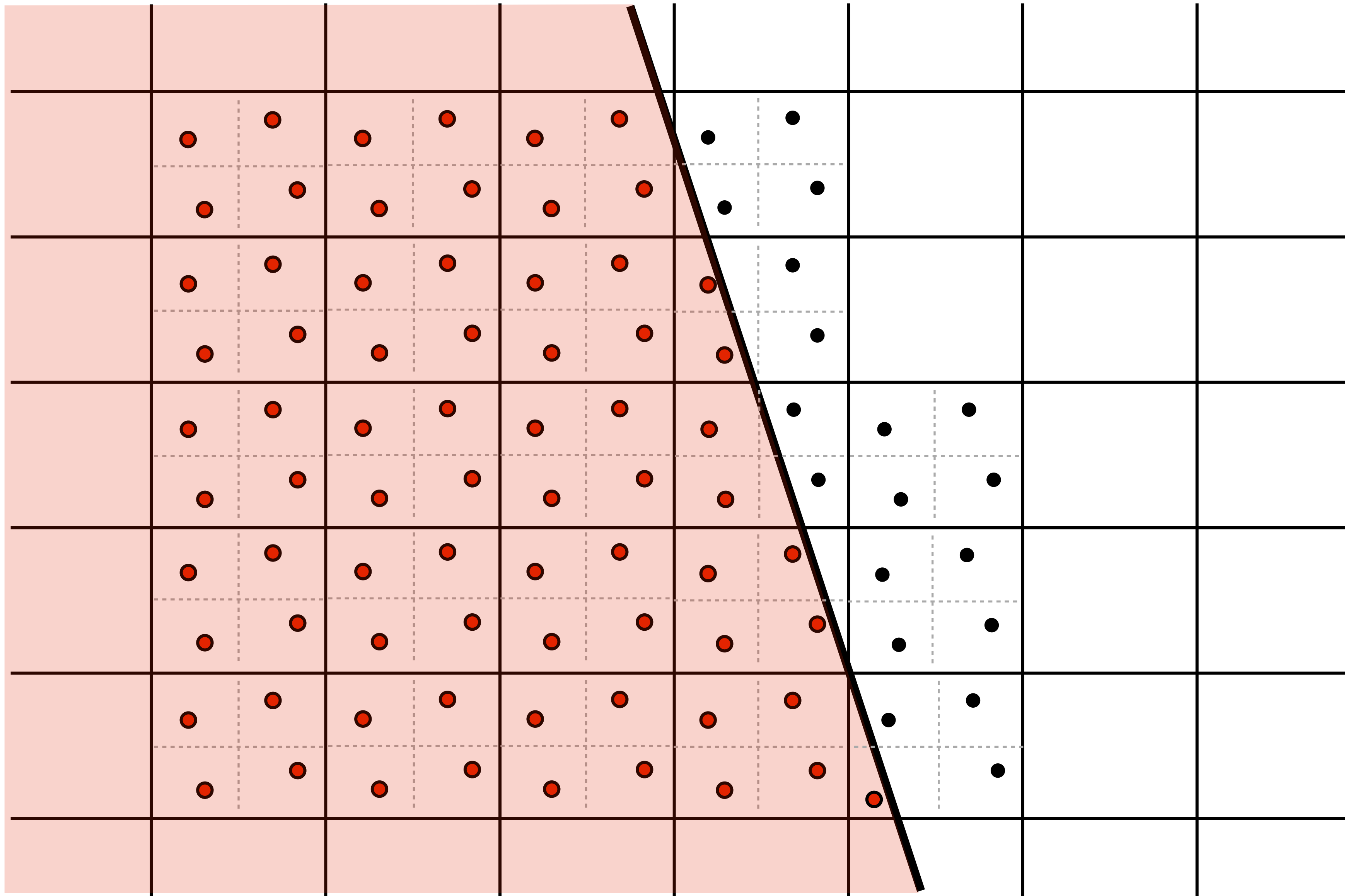
- **Undersampling high frequency signal results in aliasing**
 - **“Jaggies” in a single image**
 - **“Roping” or “shimmering” in an animation**

Supersampling: increase rate of sampling to more accurately reconstruct high frequencies in triangle coverage signal (high frequencies exist because of triangle edges)



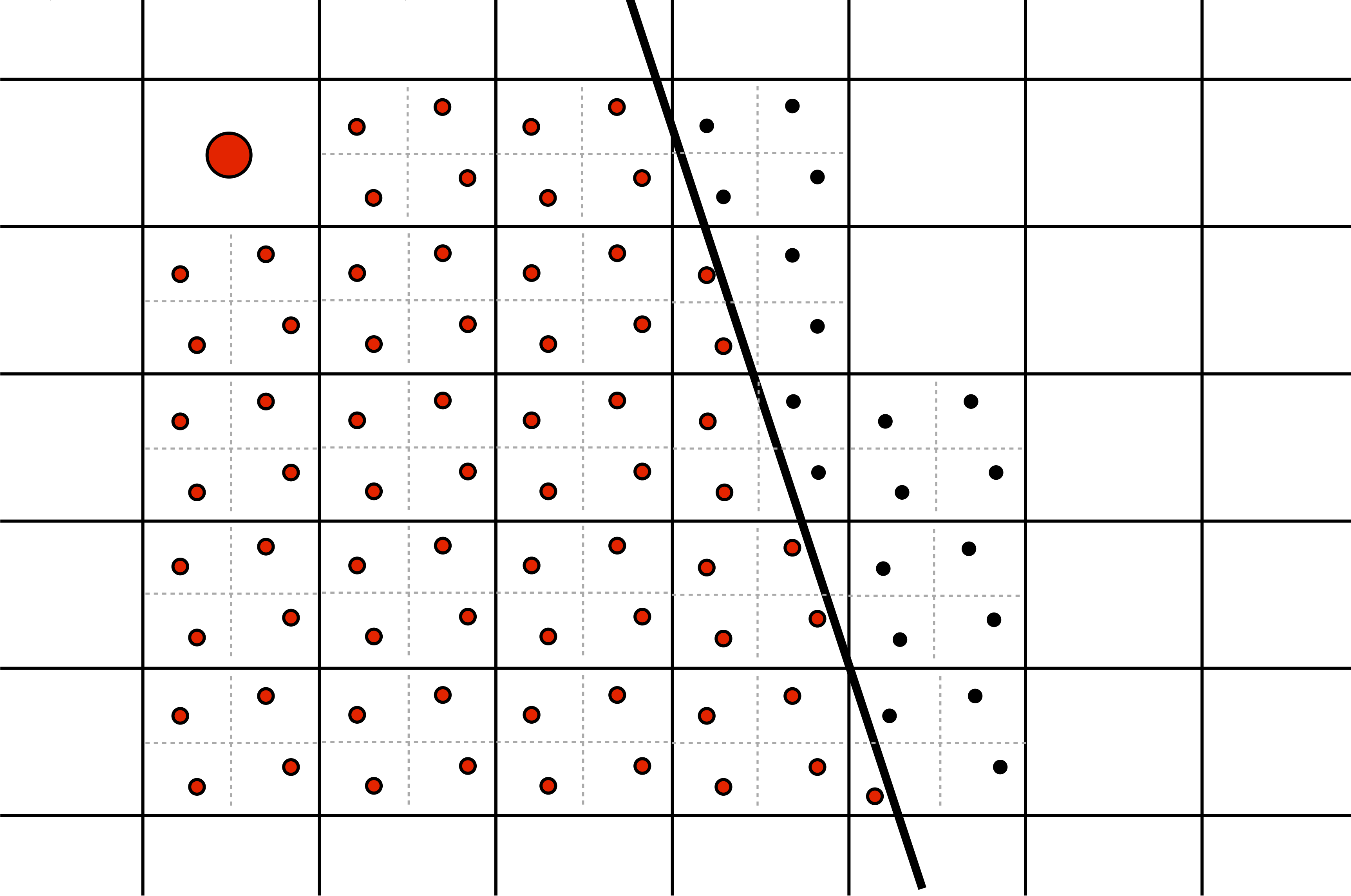
Supersampling

Example: stratified sampling using
four samples per pixel

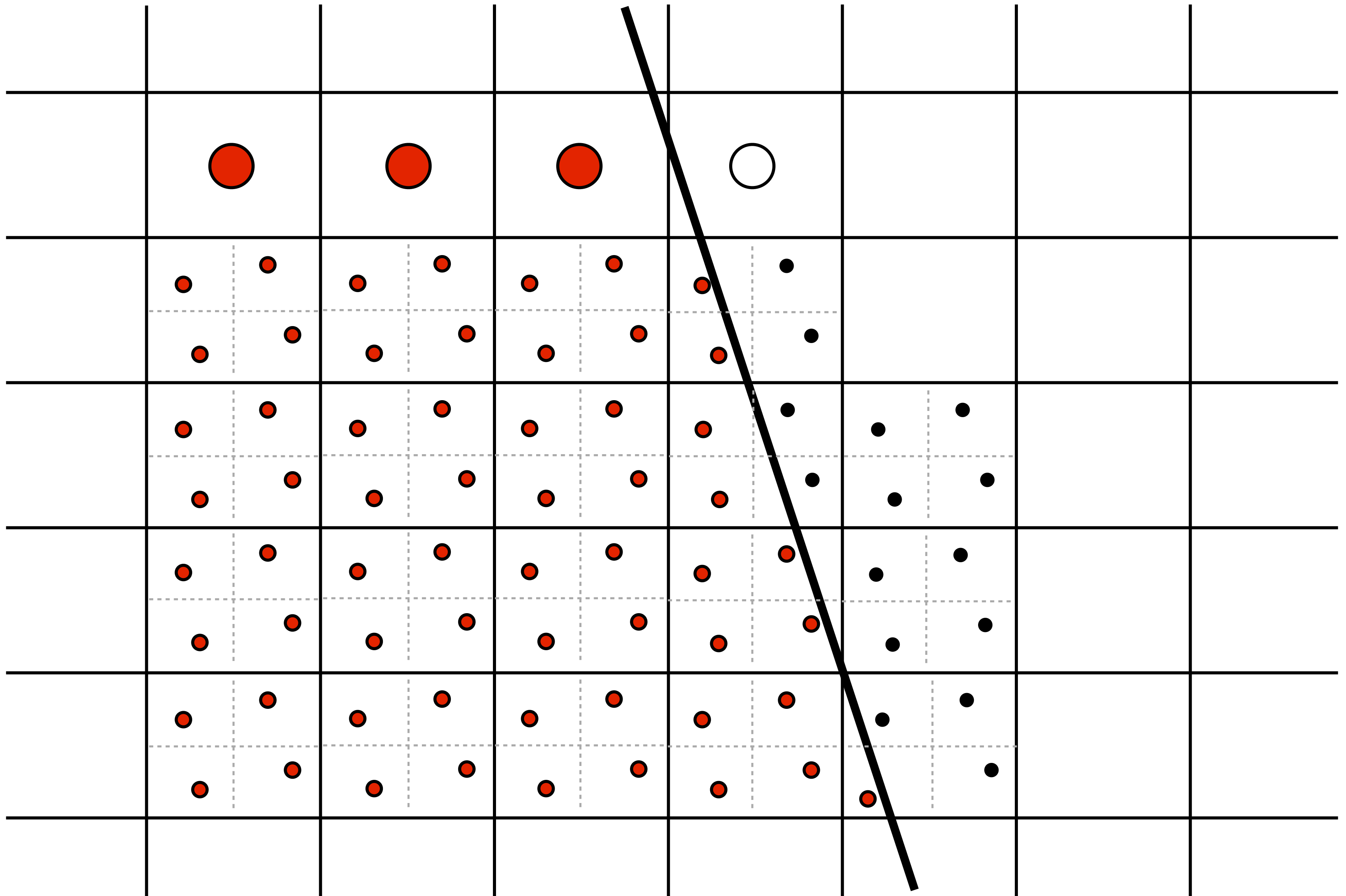


Resample to display's pixel rate (using box filter)

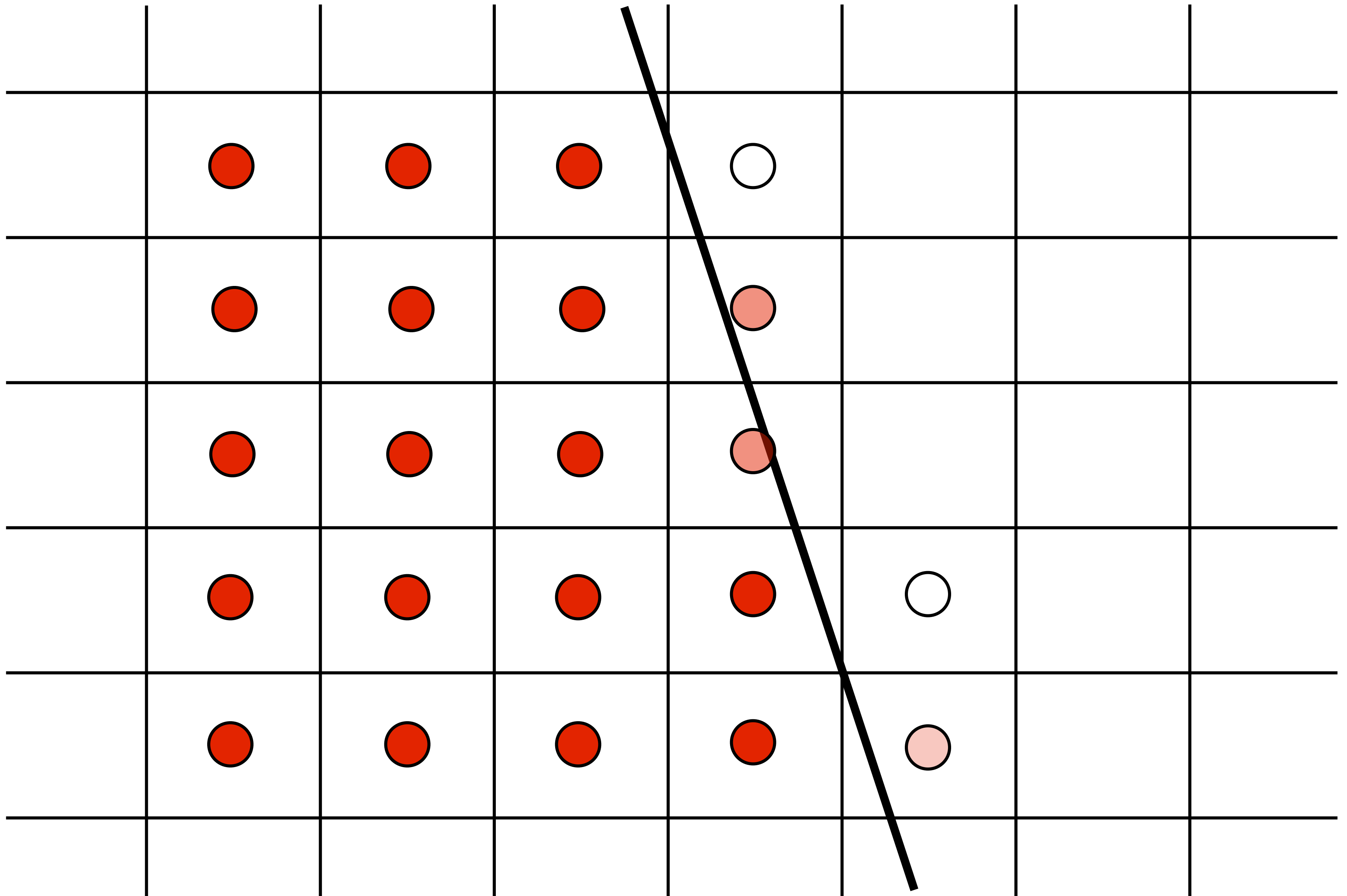
(Why? Because a screen displays one sample value per screen pixel... that's the definition of a pixel)



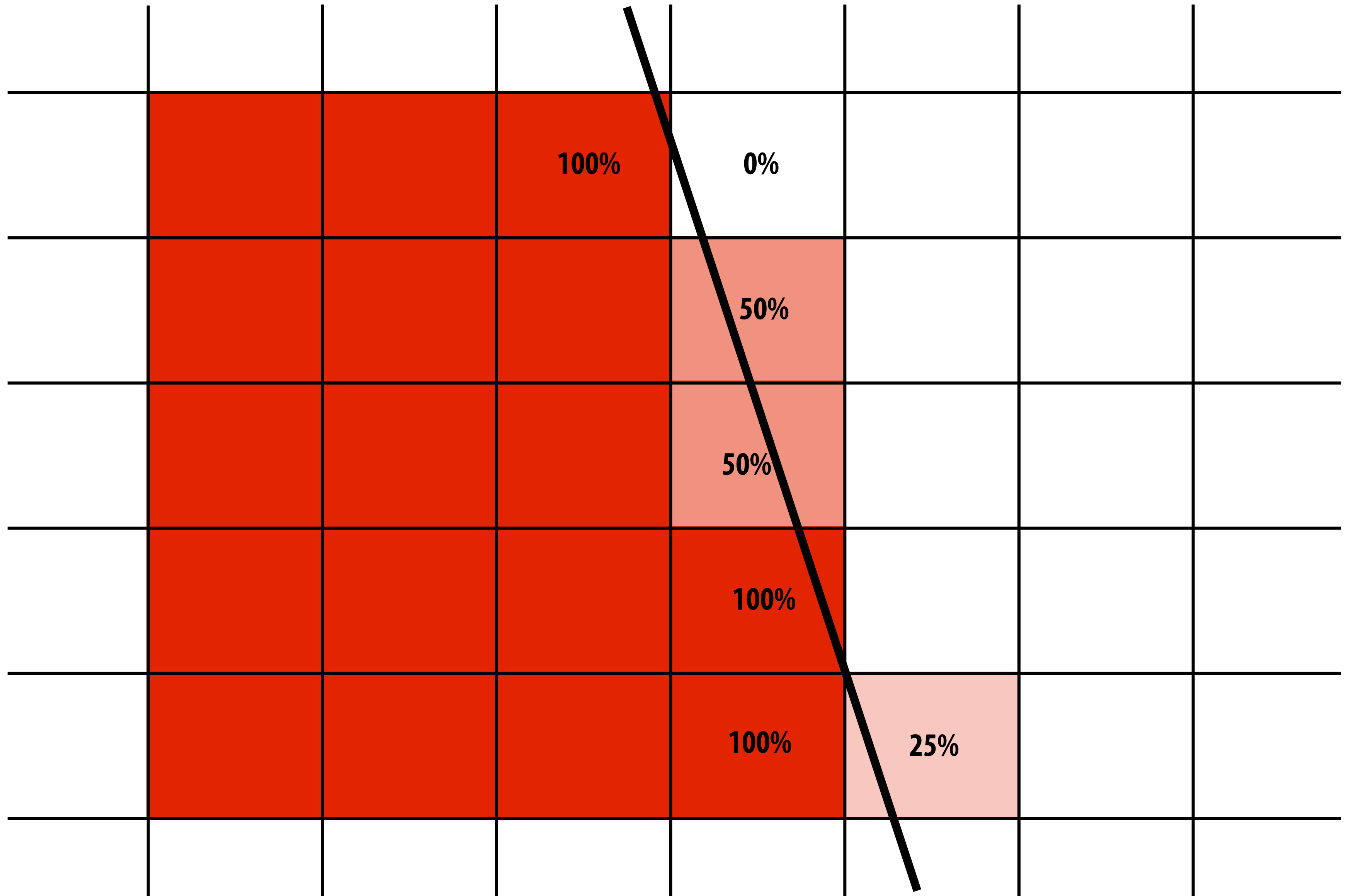
Resample to display's pixel rate (using box filter)



Resample to display's pixel rate (using box filter)



Displayed result (note anti-aliased edges)



Sampling coverage

- **What we really want is for actual displayed intensity of a region of the physical screen to closely approximate the exact intensity of that region as measured by the scene's virtual camera.**
- **So we want to produce values to send to display that approximate the integral of scene radiance for the region illuminated by a display pixel (supersampling is used to estimate this integral)**

Modes of fragment generation

- **Supersampling: generate one fragment for each covered sample**
- **Multi-sampling: generate one fragment per pixel if any sample point within the pixel is covered**
- **Today, let's assume that the number of samples per pixel is one.
(thus, both of the above schemes are equivalent)**

Point-in-triangle test

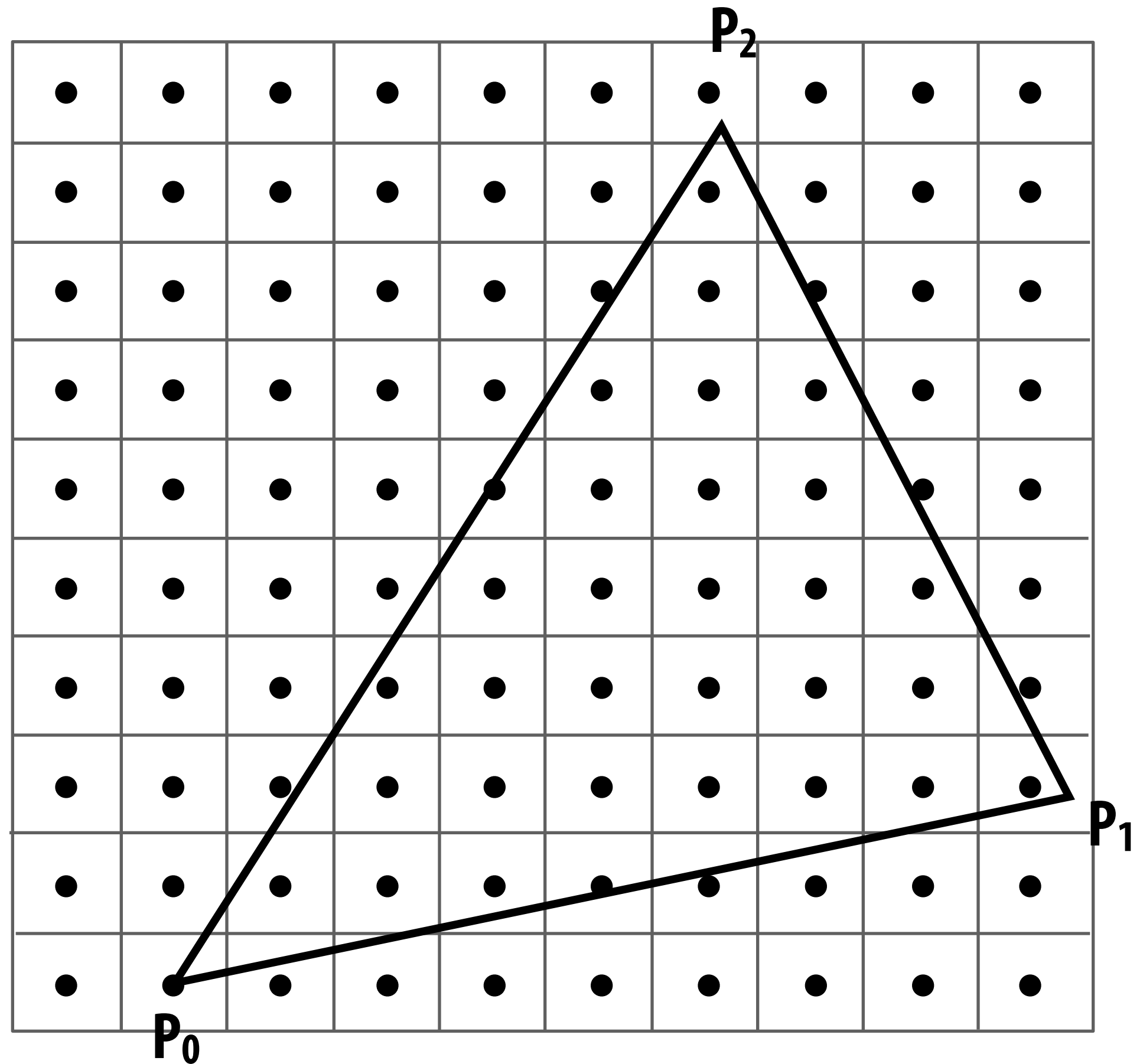
$$P_i = (x_i/w_i, y_i/w_i, z_i/w_i) = (X_i, Y_i, Z_i)$$

$$dX_i = X_{i+1} - X_i$$

$$dY_i = Y_{i+1} - Y_i$$

$$\begin{aligned} E_i(x, y) &= (x - X_i) dY_i - (y - Y_i) dX_i \\ &= A_i x + B_i y + C_i \end{aligned}$$

$$\begin{aligned} E_i(x, y) = 0 &: \text{point on edge} \\ > 0 &: \text{outside edge} \\ < 0 &: \text{inside edge} \end{aligned}$$



Point-in-triangle test

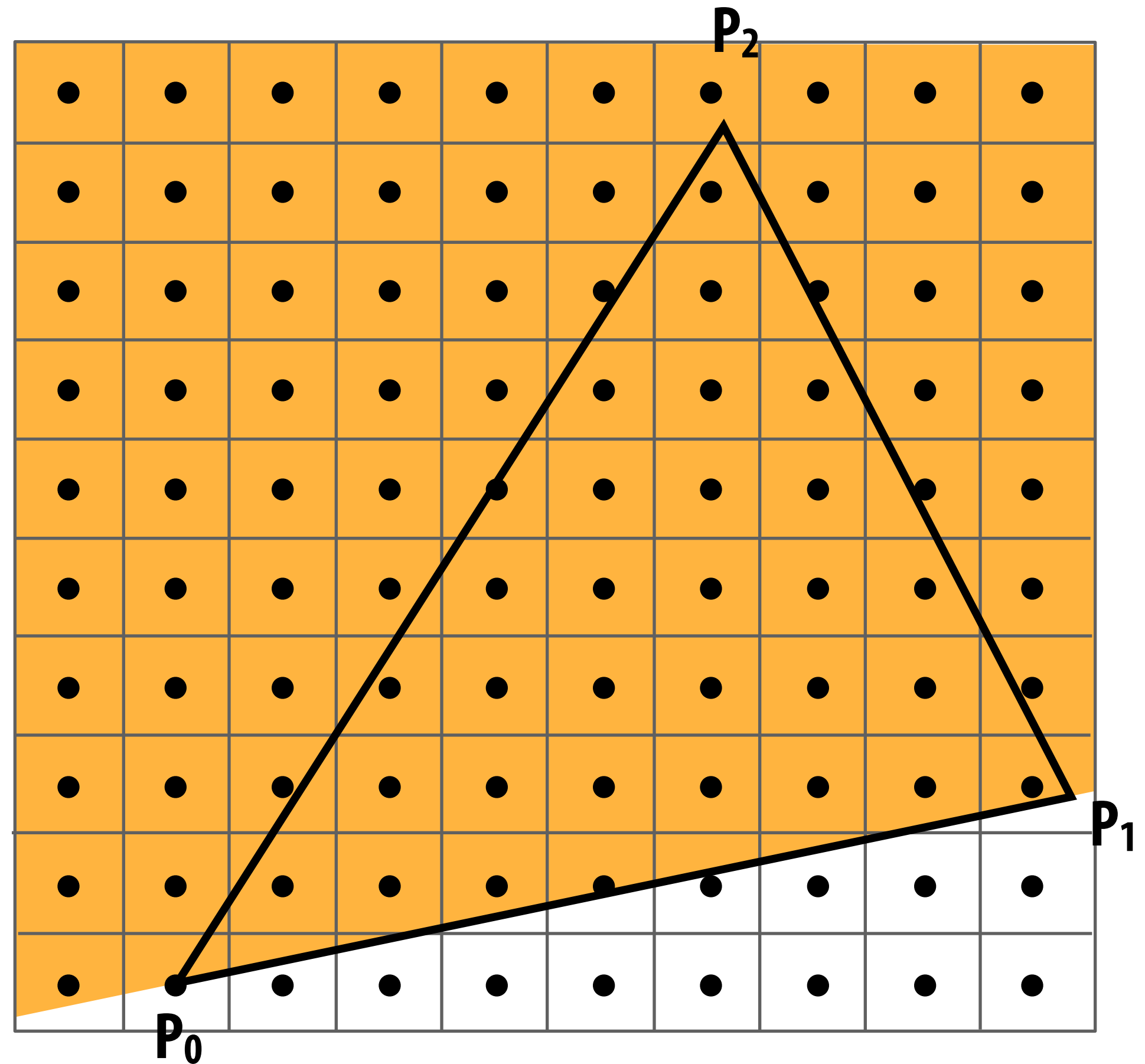
$$P_i = (x_i/w_i, y_i/w_i, z_i/w_i) = (X_i, Y_i, Z_i)$$

$$dX_i = X_{i+1} - X_i$$

$$dY_i = Y_{i+1} - Y_i$$

$$\begin{aligned} E_i(x, y) &= (x - X_i) dY_i - (y - Y_i) dX_i \\ &= A_i x + B_i y + C_i \end{aligned}$$

$$\begin{aligned} E_i(x, y) = 0 &: \text{point on edge} \\ > 0 &: \text{outside edge} \\ < 0 &: \text{inside edge} \end{aligned}$$



Point-in-triangle test

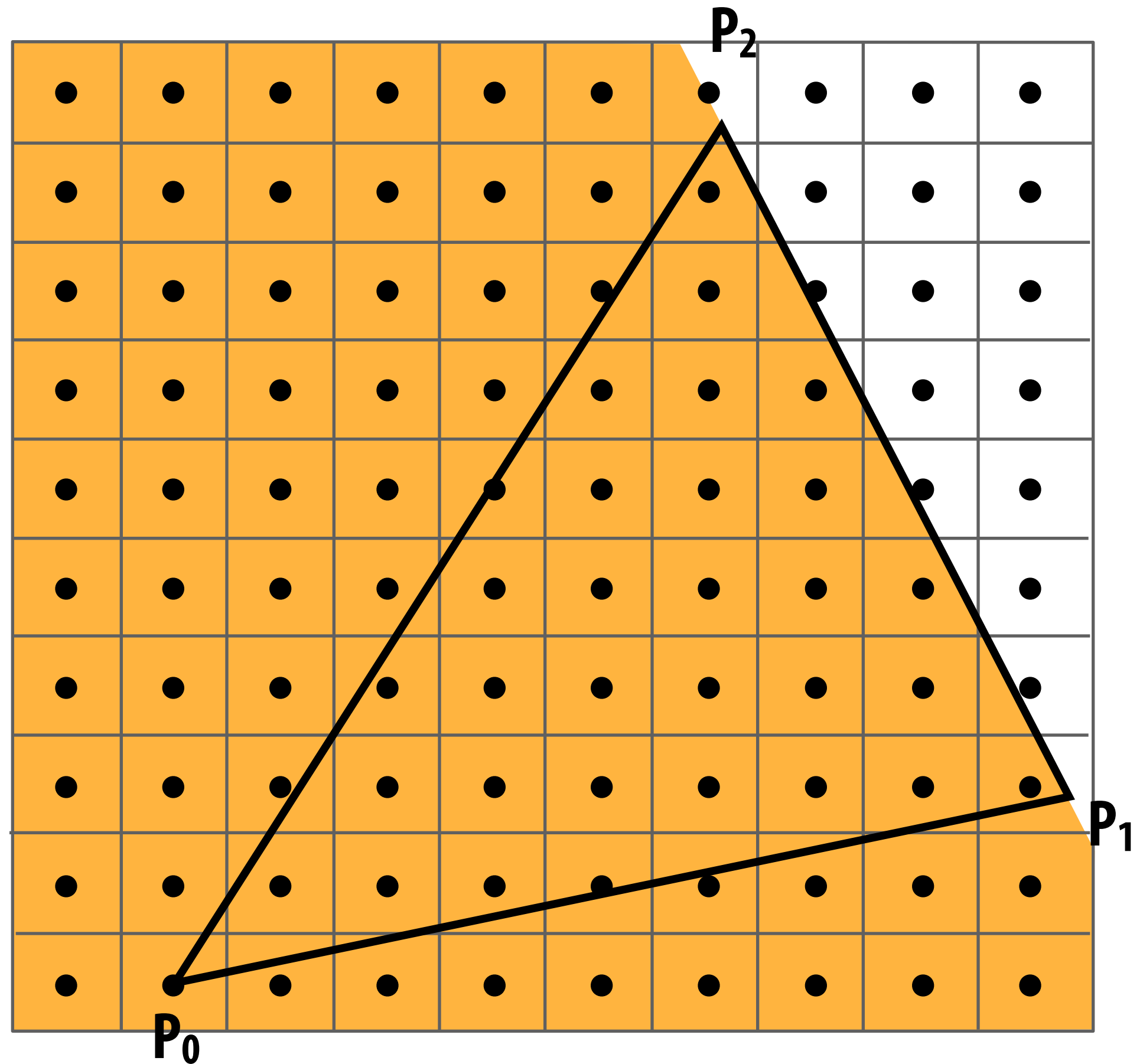
$$P_i = (x_i/w_i, y_i/w_i, z_i/w_i) = (X_i, Y_i, Z_i)$$

$$dX_i = X_{i+1} - X_i$$

$$dY_i = Y_{i+1} - Y_i$$

$$\begin{aligned} E_i(x, y) &= (x - X_i) dY_i - (y - Y_i) dX_i \\ &= A_i x + B_i y + C_i \end{aligned}$$

$$\begin{aligned} E_i(x, y) &= 0 : \text{point on edge} \\ &> 0 : \text{outside edge} \\ &< 0 : \text{inside edge} \end{aligned}$$



Point-in-triangle test

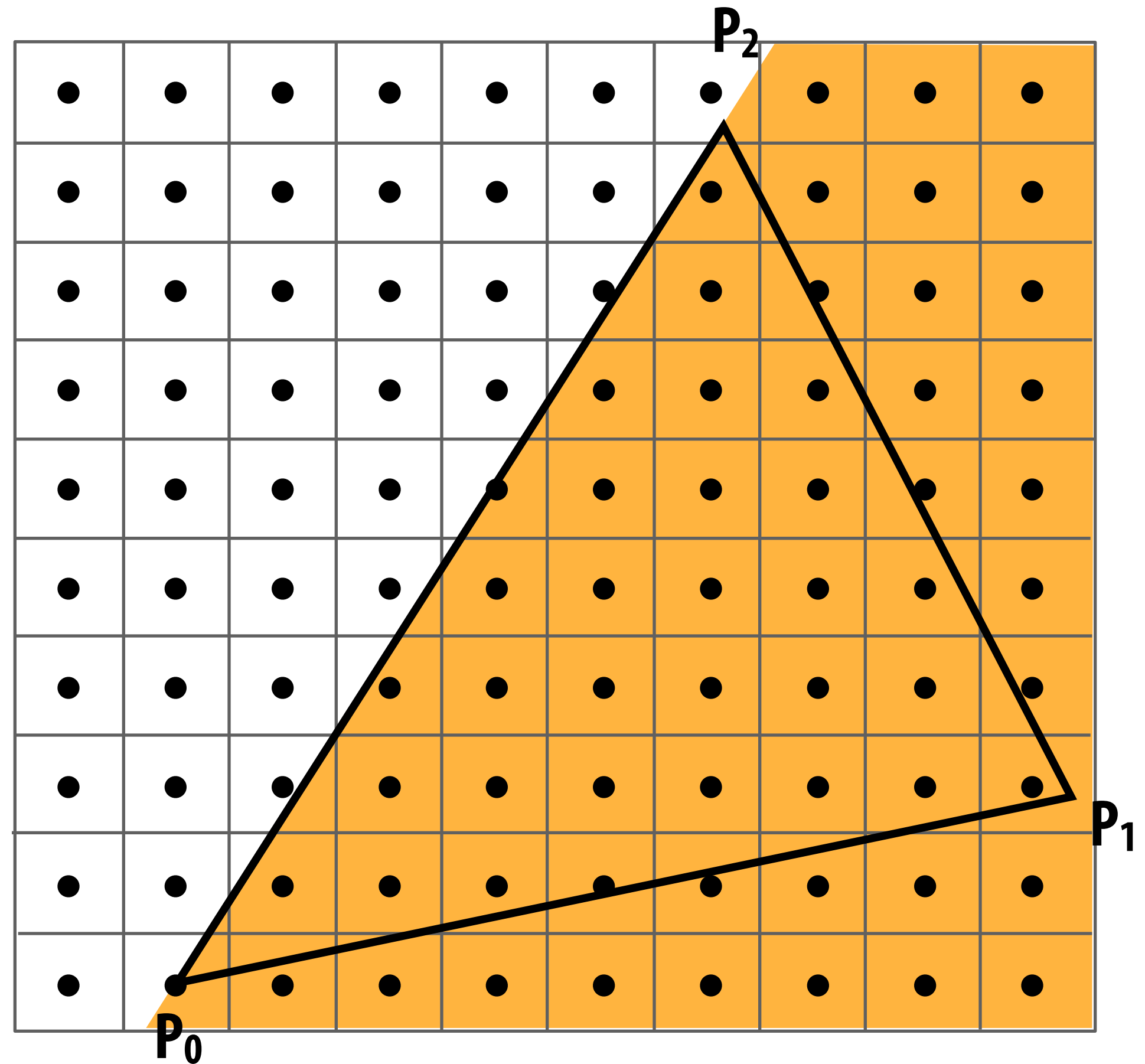
$$P_i = (x_i/w_i, y_i/w_i, z_i/w_i) = (X_i, Y_i, Z_i)$$

$$dX_i = X_{i+1} - X_i$$

$$dY_i = Y_{i+1} - Y_i$$

$$\begin{aligned} E_i(x, y) &= (x - X_i) dY_i - (y - Y_i) dX_i \\ &= A_i x + B_i y + C_i \end{aligned}$$

$$\begin{aligned} E_i(x, y) &= 0 : \text{point on edge} \\ &> 0 : \text{outside edge} \\ &< 0 : \text{inside edge} \end{aligned}$$

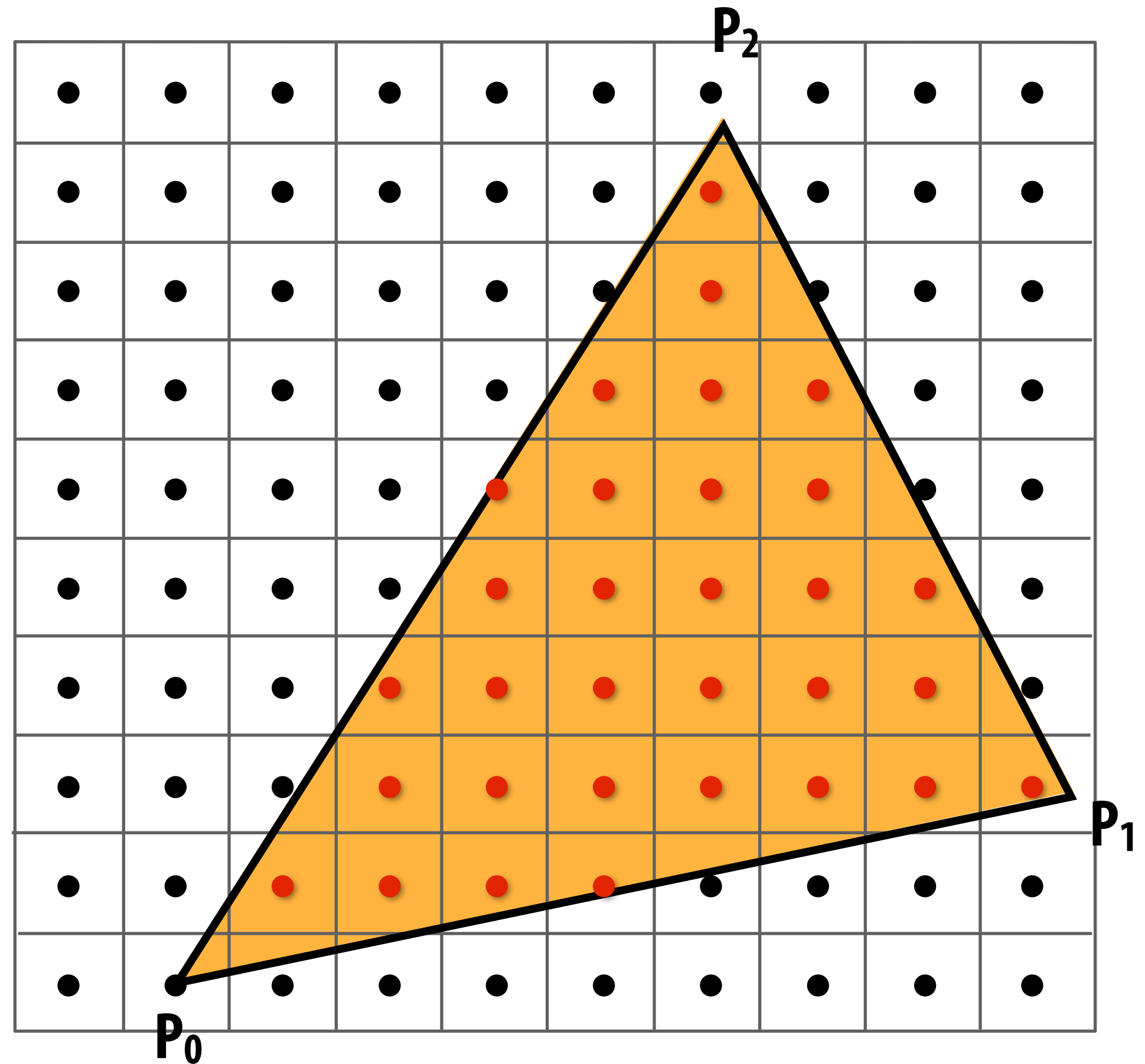


Point-in-triangle test

Sample point $s = (sx, sy)$ is inside the triangle if it is inside all three edges.

$inside(sx, sy) =$
 $E_0(sx, sy) < 0 \ \&\&$
 $E_1(sx, sy) < 0 \ \&\&$
 $E_2(sx, sy) < 0;$

Note: actual implementation of $inside(sx, sy)$ involves \leq checks based on the pipeline rasterizer's edge rules.



Sample points inside triangle are highlighted red.

Incremental triangle traversal

$$P_i = (x_i/w_i, y_i/w_i, z_i/w_i) = (X_i, Y_i, Z_i)$$

$$dX_i = X_{i+1} - X_i$$

$$dY_i = Y_{i+1} - Y_i$$

$$\begin{aligned} E_i(x, y) &= (x - X_i) dY_i - (y - Y_i) dX_i \\ &= A_i x + B_i y + C_i \end{aligned}$$

$E_i(x, y) = 0$: point on edge
 > 0 : outside edge
 < 0 : inside edge

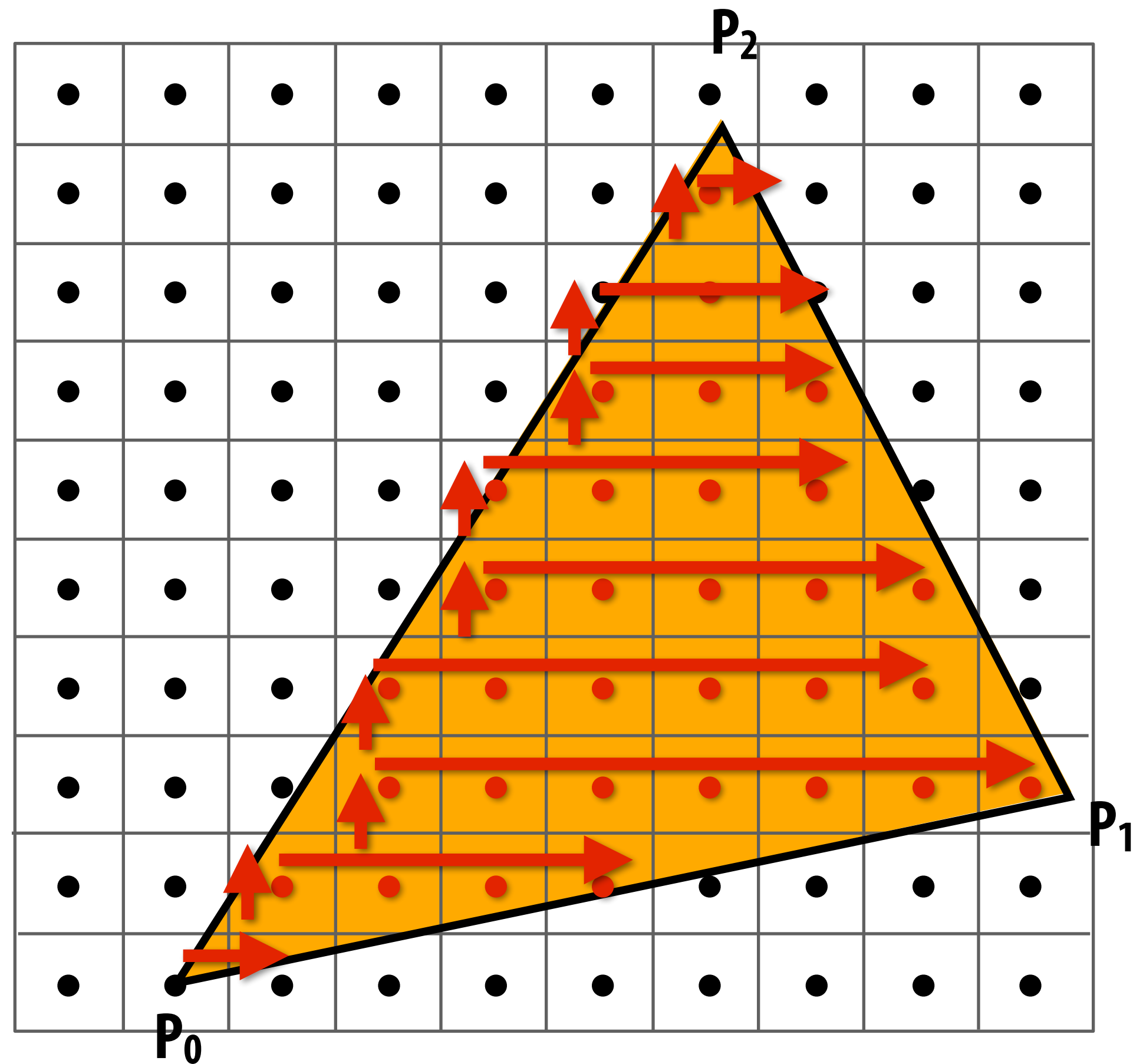
Note incremental update:

$$dE_i(x+1, y) = E_i(x, y) + dY_i = E_i(x, y) + A_i$$

$$dE_i(x, y+1) = E_i(x, y) + dX_i = E_i(x, y) + B_i$$

Incremental update saves computation:
One addition per edge, per sample test

Note: many traversal orders are possible: backtrack, zig-zag, Hilbert/Morton curves (locality maximizing)



Modern hierarchical traversal

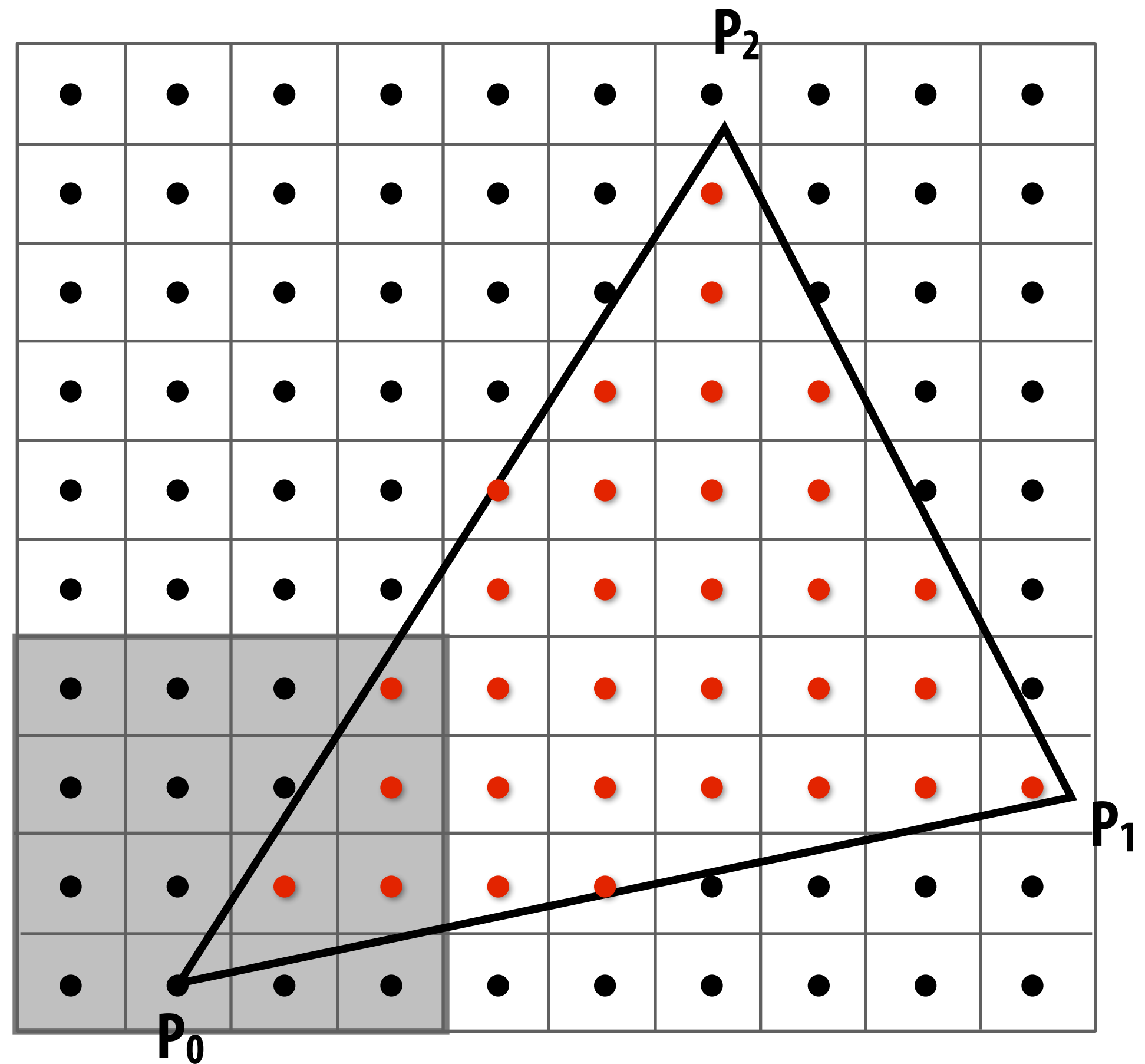
Traverse triangle in blocks

**Test all samples in block against triangle in parallel
(e.g., data-parallel hardware implementation)**

Can be implemented as multi-level hierarchy.

Advantages:

- **Simplicity of wide parallel execution overcomes cost of extra point-in-triangle tests (recall: most triangles cover many samples, especially when super-sampling coverage)**
- **Can skip sample testing work: entire block not in triangle ("early out"), entire block entirely within triangle ("early in")**
- **Entire tile may be discarded due to occlusion (see occlusion cull later in this lecture)**



Sampling triangle attributes

Computing attributes

- **How are fragment attributes (color, normal, texcoords) computed?**
 - Point sample attributes as well (evaluate attributes at sample point)
 - Fragment contains $A(x,y)$ for all attributes (needed for fragment shading)

Computing a plane equation for an attribute:

Let A_0, A_1, A_2 be attribute values at the three triangle vertices

Let projected screen-space positions of vertices be $(X_0, Y_0), (X_1, Y_1), (X_2, Y_2)$

Linear interpolation of vertex attributes, so $A(x,y) = ax + by + c$ (attribute plane equation)

$$A_0 = aX_0 + bY_0 + c$$

$$A_1 = aX_1 + bY_1 + c$$

$$A_2 = aX_2 + bY_2 + c$$

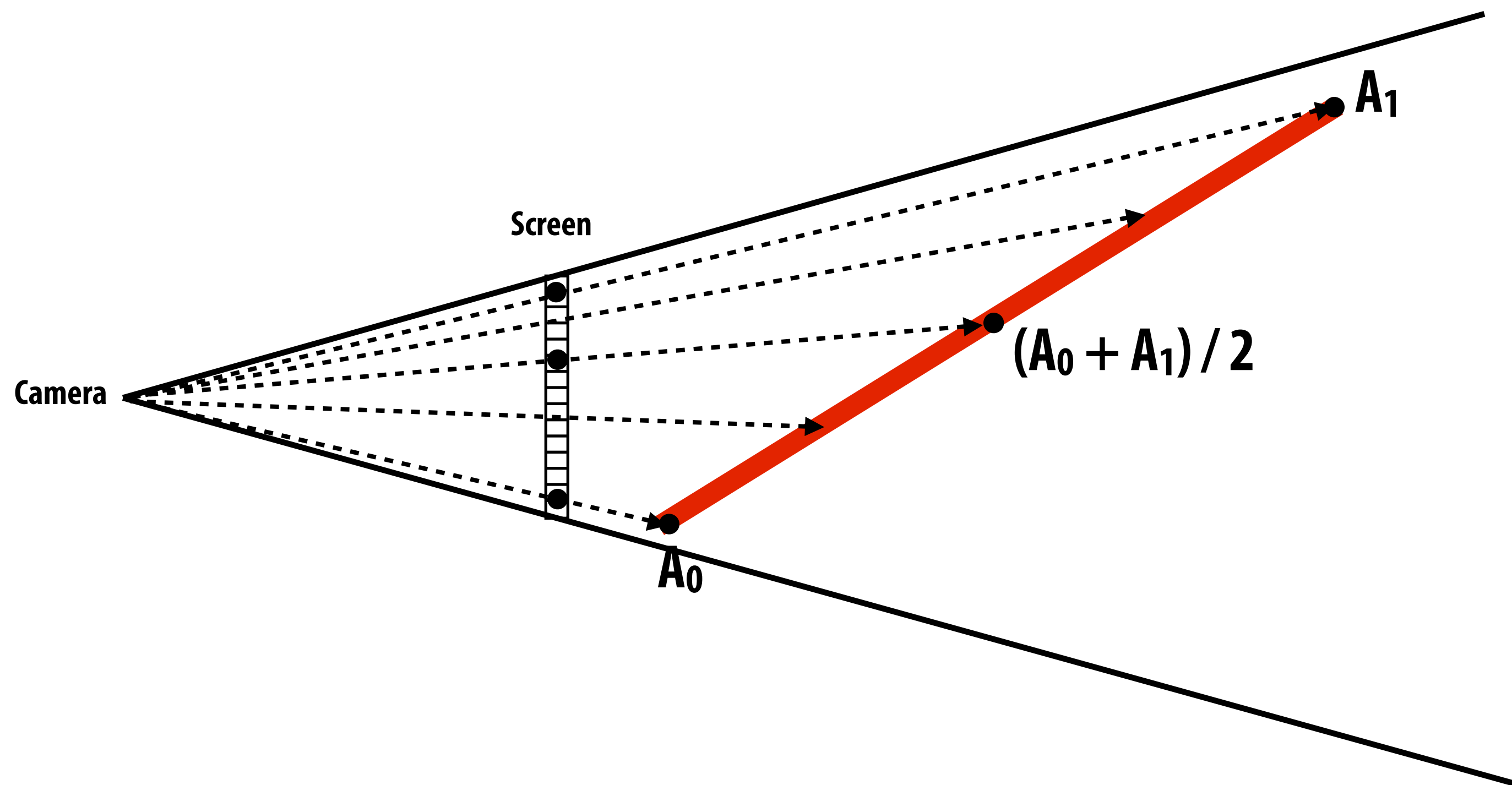
3 equations, 3 unknowns. Solve for a,b,c **

** Discard zero-area triangles before getting here (recall we computed area in back-face culling)

Perspective-correct interpolation

Due to projection, interpolation is not linear in screen XY coordinates

Attribute values should be interpolated linearly on triangle in 3D object space.



Perspective-correct interpolation

Assume triangle attribute varies linearly across the triangle:

Attribute's value at 3D (non-homogeneous) point $P=(x, y, z)$: $A = ax + by + cz$

Rewrite attribute plane equation in terms of 2D homogeneous coordinates:

Project P , get 2D homogeneous representation $(x_{2dh}, y_{2dh}, w) = (x, y, z)$

So... $A = ax_{2dh} + by_{2dh} + cw$

$$\frac{A}{w} = a \left(\frac{x_{2dh}}{w} \right) + b \left(\frac{y_{2dh}}{w} \right) + c$$

$$\frac{A}{w} = aX + bY + c \quad \text{Where } (X, Y) \text{ are projected screen 2D coordinates (after homogeneous divide)}$$

So ... $\frac{A}{w}$ is linear in projected screen 2D coordinates (X, Y)

Perspective-correct interpolation

Attribute values vary linearly across triangle in 3D, but not linear in projected screen XY

Projected attribute values (A/w) are linear in screen XY!

For each generated fragment:

Evaluate $1/w(x,y)$ (from precomputed plane equation for $1/w$)

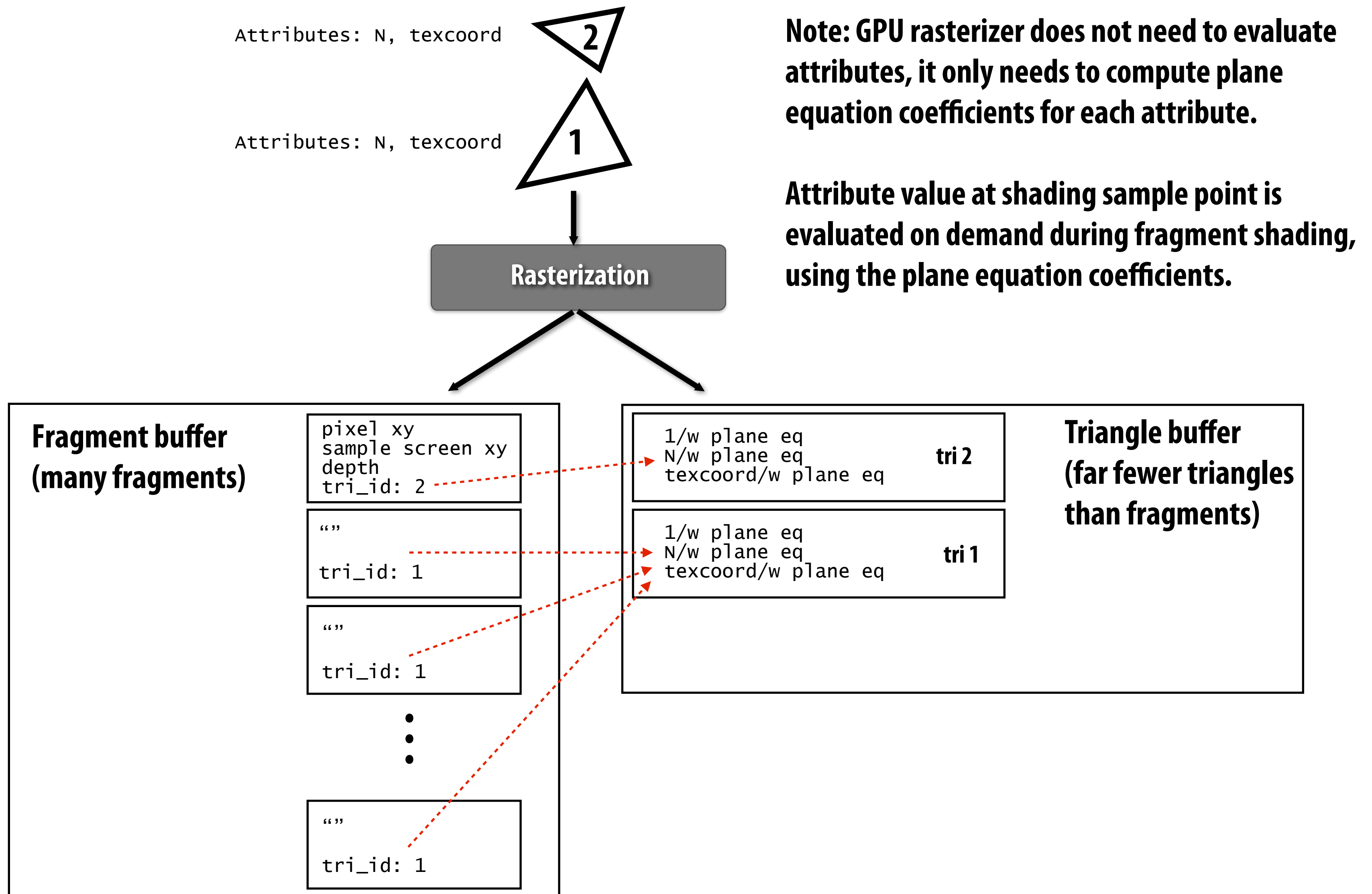
Reciprocate $1/w(x,y)$ to get $w(x,y)$

For each triangle attribute:

Evaluate $A/w(x,y)$ (from precomputed plane equation for A/w)

Multiply $A/w(x,y)$ by $w(x,y)$ to get $A(x,y)$

GPUs store attribute plane equations separately from individual fragments (effectively a form of fragment compression)



Modern GPU rasterization

■ Triangle setup:

- Transform clip space vertex positions to screen space
- Convert vertex positions to fixed point (Direct3D requires 8 bits of subpixel precision**)
- Compute triangle edge equations
- Compute plane equations for all vertex attributes, $1/w$, and Z

■ Traverse triangle in blocks:

- May attempt to trivially accept/reject block using edge tests on block corners
- Identify covered samples using edge tests (wide data parallelism in implementation)
- Generate and emit fragments based on coverage (also emit per-triangle data as necessary)
- Block granularity order is also important for... (topics of future lectures)
 - Shading derivatives, maximizing data locality (cache locality/compression), maximizing control locality (avoiding SIMD divergence during fragment shading)

**** Note 1: limited precision can be a good thing: really acute triangles snap to 0 area and get discarded**

**** Note 2: limited precision can be a bad thing: precision limits in (x,y) can limit precision in Z (see Akeley and Su, 2006)**

Occlusion

Depth buffer for occlusion

- **Depth buffer stores depth of scene at each coverage sample point**
 - Stored per sample, not per pixel!
- **Triangles are planar**
 - Each triangle has exactly one depth at each sample point * ✓
(so triangle order is a well-defined ordering of fragments at each sample point)
- **Occlusion check using Z-buffer algorithm**
 - Constant-time occlusion test per fragment ✓
 - Constant space per coverage sample ✓
 - Constant space per depth buffer (overall) ✓

* Assumes edge-on triangles have been discarded due to zero area

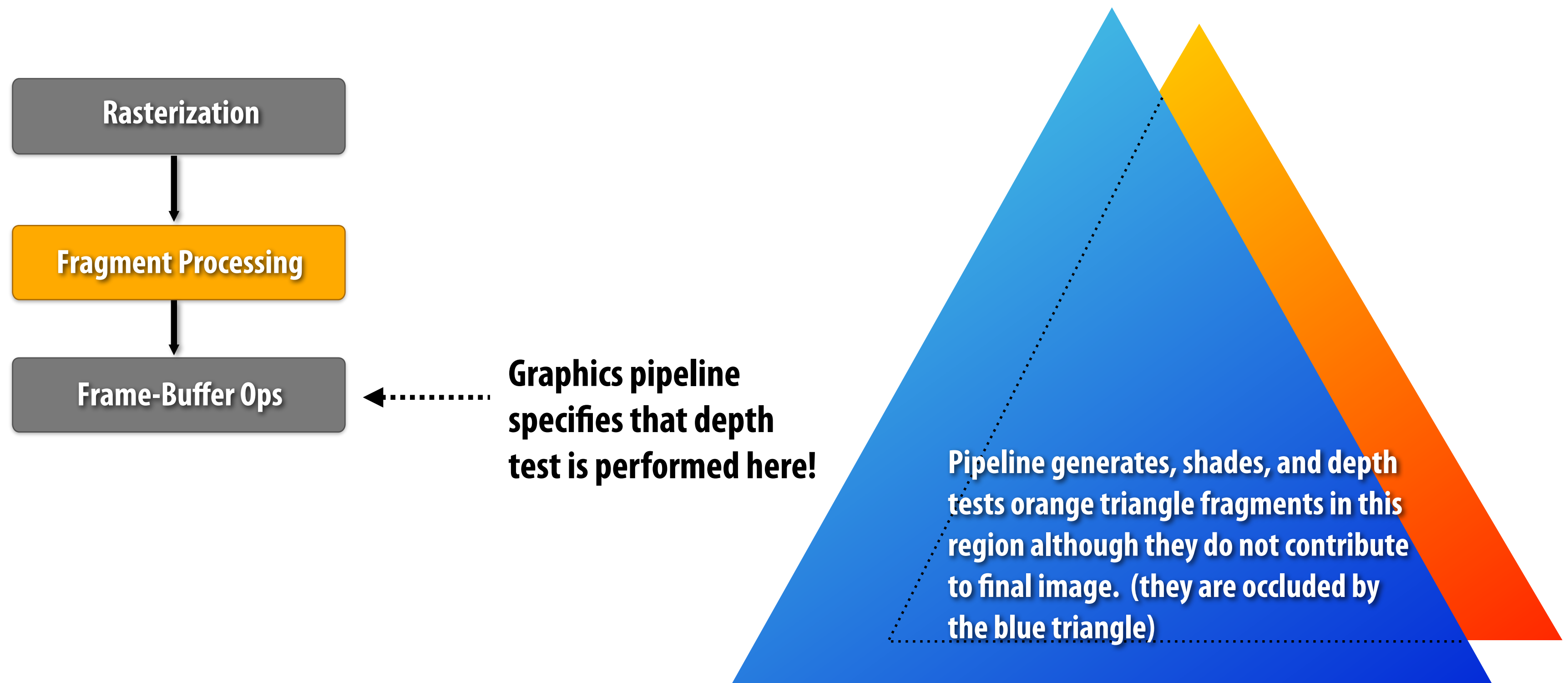
Depth buffer for occlusion

- **Z-buffer algorithm has high bandwidth requirements (particularly when super-sampling triangle coverage)**
 - Number of Z-buffer reads/writes for a frame depends on:
 - Depth complexity of the scene
 - The order triangles are provided to the graphics pipeline
(if depth test fails, don't write to depth buffer or rgba)
- **Bandwidth estimate:**
 - $60 \text{ Hz} \times 2 \text{ MPixel image} \times \text{avg. depth complexity } 4 \text{ (assume: replace 50\% of time)} \times 32\text{-bit Z} = 2.8 \text{ GB/s}$
 - If super-sampling at 4 times per pixel, multiply by 4
 - Consider five shadow maps per frame (1 MPixel, not super-sampled): additional 8.6 GB/s
 - Note: this is just depth accesses. It does not include color-buffer bandwidth
- **Modern GPUs implement caching and lossless compression of both color and depth buffers to reduce bandwidth (later this lecture)**

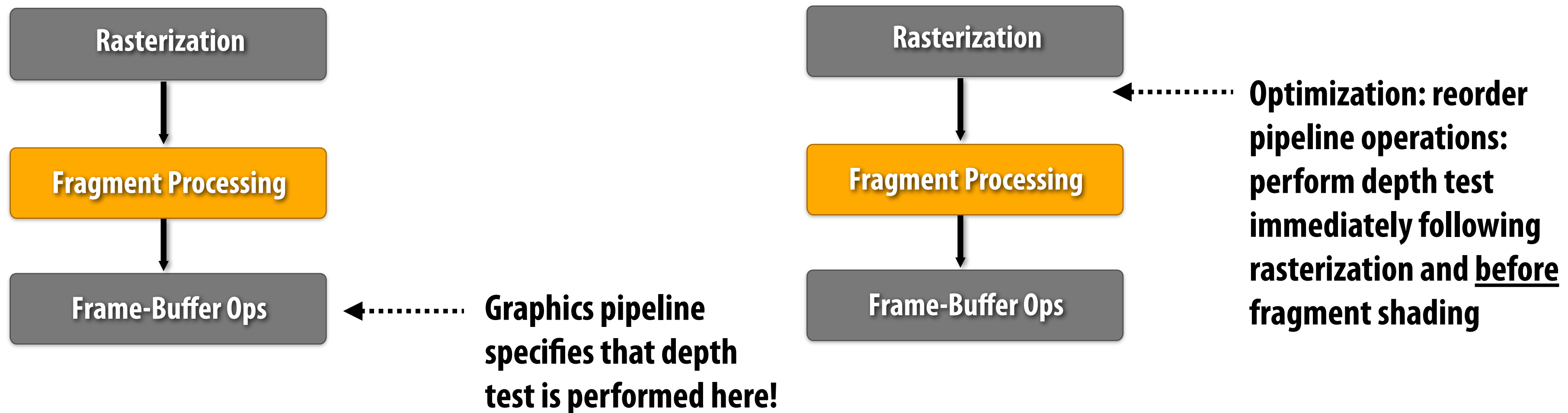
Early occlusion culling

Early occlusion-culling (“early Z”)

Idea: discard fragments that will not contribute to image as quickly as possible in the pipeline



Early occlusion-culling (“early Z”)



A GPU implementation detail: not reflected in the graphics pipeline abstraction

Key assumption: occlusion results do not depend on fragment shading

- Example operations that prevent use of this early Z optimization: enabling alpha test, fragment shader modifies fragment's Z value

Note: early Z only provides benefit if closer triangle is rendered by application first!

(application developers are encouraged to submit geometry in as close to front-to-back order as possible)

Summary: early occlusion culling

- **Key observation: can reorder pipeline operations without impacting correctness: perform depth test prior to fragment shading**
- **Benefit: reduces fragment processing work**
 - Effectiveness of optimization is dependent on triangle ordering
 - Ideal geometry submission order: front-to-back order
- **Does not reduce amount of bandwidth used to perform depth tests**
 - The same depth-buffer reads and writes still occur (they just occur before fragment shading)
- **Implementation-specific optimization, but programmers know it is there**
 - Commonly used two-pass technique: **rendering with a “Z-prepass”**
 - Pass 1: render all scene geometry, with fragment shading and color buffer writes disabled (Put the depth buffer in its end-of-frame state)
 - Pass 2: re-render scene with shading enabled and with depth-test predicate: less than-or-equal
 - Overhead: must process and rasterizer scene geometry twice
 - Benefit: minimizes expensive fragment shading work by only shading visible fragments

Hierarchical early occlusion culling: “hi-Z”

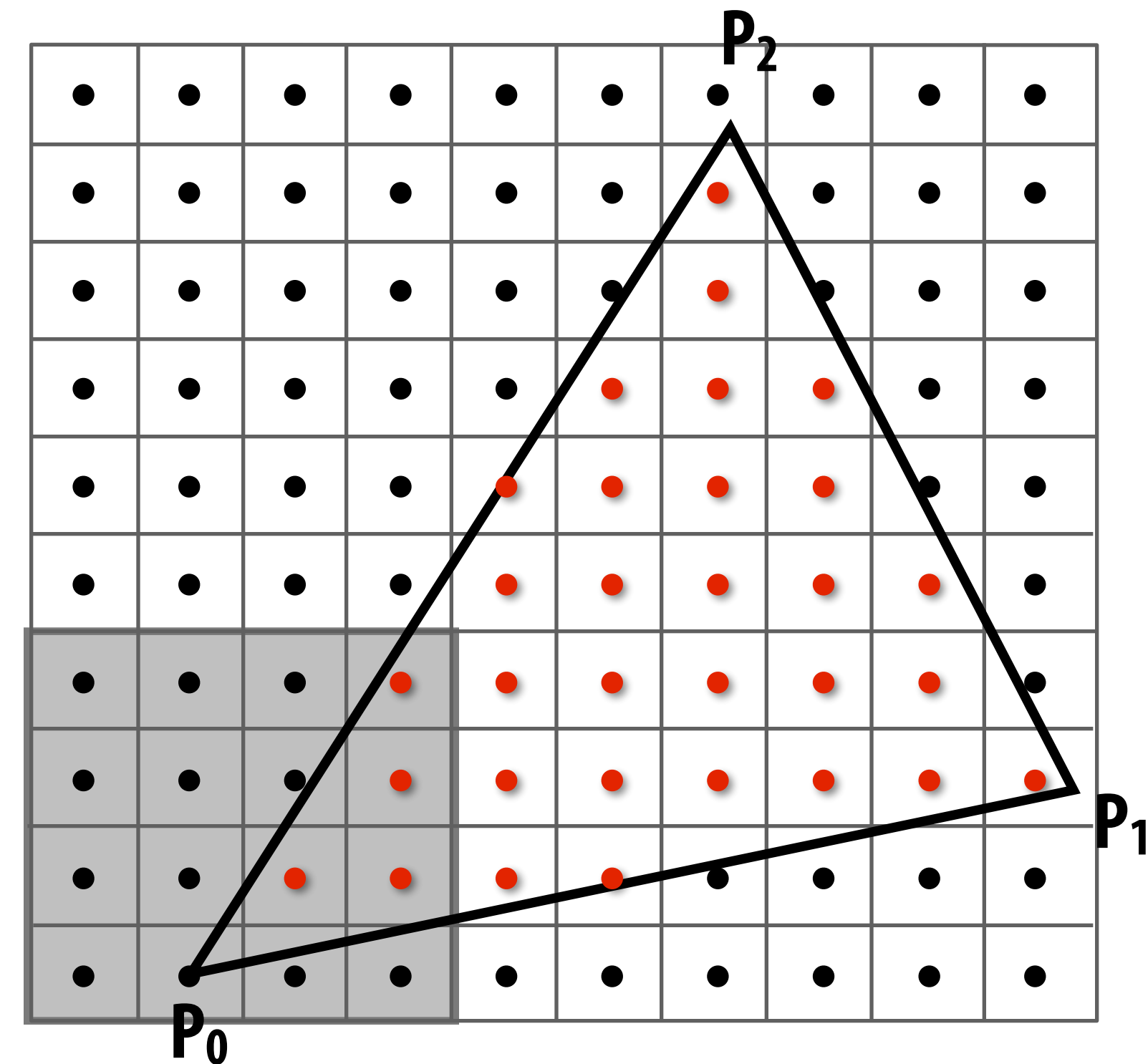
Recall hierarchical traversal during rasterization

Z-Max culling:

For each screen tile, compute farthest value in the depth buffer: z_{\max}

During traversal, for each tile:

1. Compute closest point on triangle in tile: tri_min (using Z plane equation)
2. If $\text{tri_min} > z_{\max}$, then triangle is completely occluded in this tile. (The depth test will fail for all samples in the tile.) Proceed to next tile without performing coverage tests for individual samples in tile.

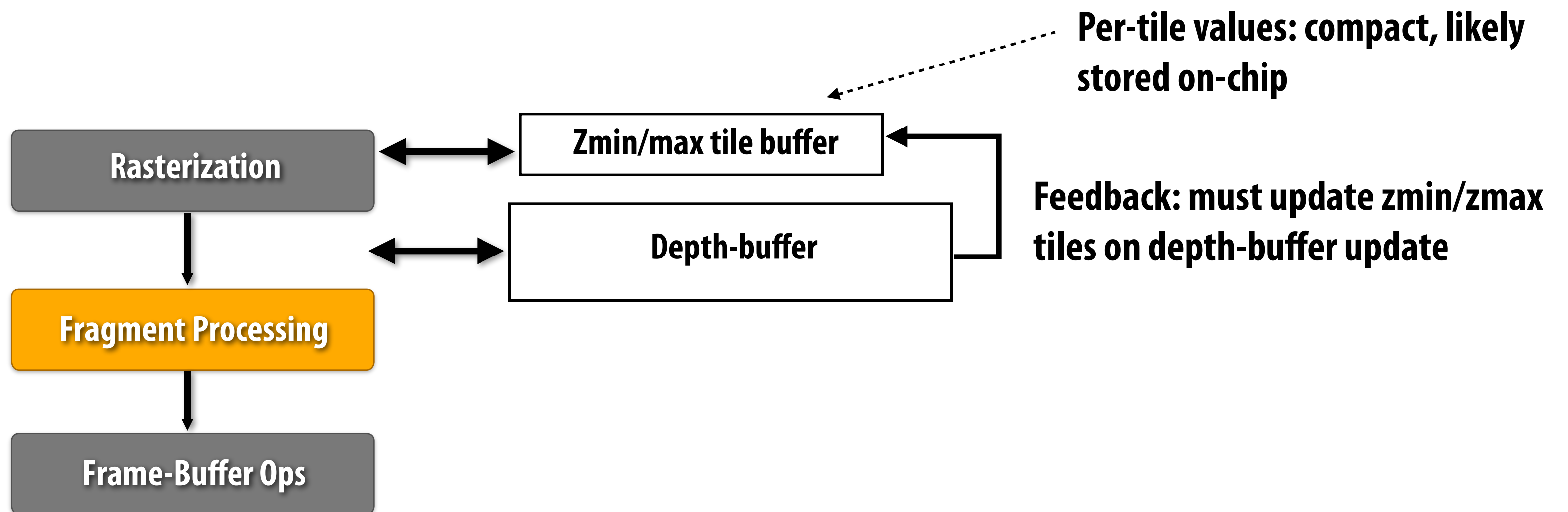


Z-min optimization:

Depth-buffer also stores z_{\min} for each tile.

If $\text{tri_max} < z_{\min}$, then all depth tests for fragments in tile will pass. (No need to perform depth test on individual fragments.)

Hierarchical Z + early Z-culling



Remember: these are GPU implementation details (common optimizations performed by most GPUs). They are invisible to the programmer and not reflected in the graphics pipeline abstraction

Summary: hierarchical Z

- **Idea: perform depth test at coarse tile granularity prior to sampling coverage**
- **ZMax culling benefits:**
 - Reduces rasterization work
 - Reduces depth-testing work (don't process individual depth samples)
 - Reduces memory bandwidth requirements (don't need to read individual depth samples)
 - Eliminates less fragment processing work than early Z (Since hierarchical Z is a conservative approximation to early X results, it will only discard a subset of the fragments early Z does)
- **ZMin benefits:**
 - Reduces depth-testing work (don't need to test individual depth samples)
 - Reduces memory bandwidth (don't need to read individual depth samples, but still must write)
- **Costs:**
 - Overhead of hierarchical tests
 - Must maintain per-tile Zmin/Zmax values
 - System complexity: must update per-tile values frequently to be effective (early Z system feeds results back to hierarchical Z system)

Fast Z clear

- **Formerly an important optimization: less important in modern GPU architectures**
 - Add “cleared” bit to tile descriptor
 - `glClear(GL_DEPTH_BUFFER)` sets these bits
 - First write to depth sample in tile unsets the “cleared” bit
- **Benefits**
 - Reduces depth-buffer write bandwidth: avoid frame-buffer write on frame-buffer clear
 - Reduces depth-buffer read bandwidth by skipping first read: if “cleared” bit for tile set, GPU can initialize tile’s contents in cache without reading data (a form of lossless compression)

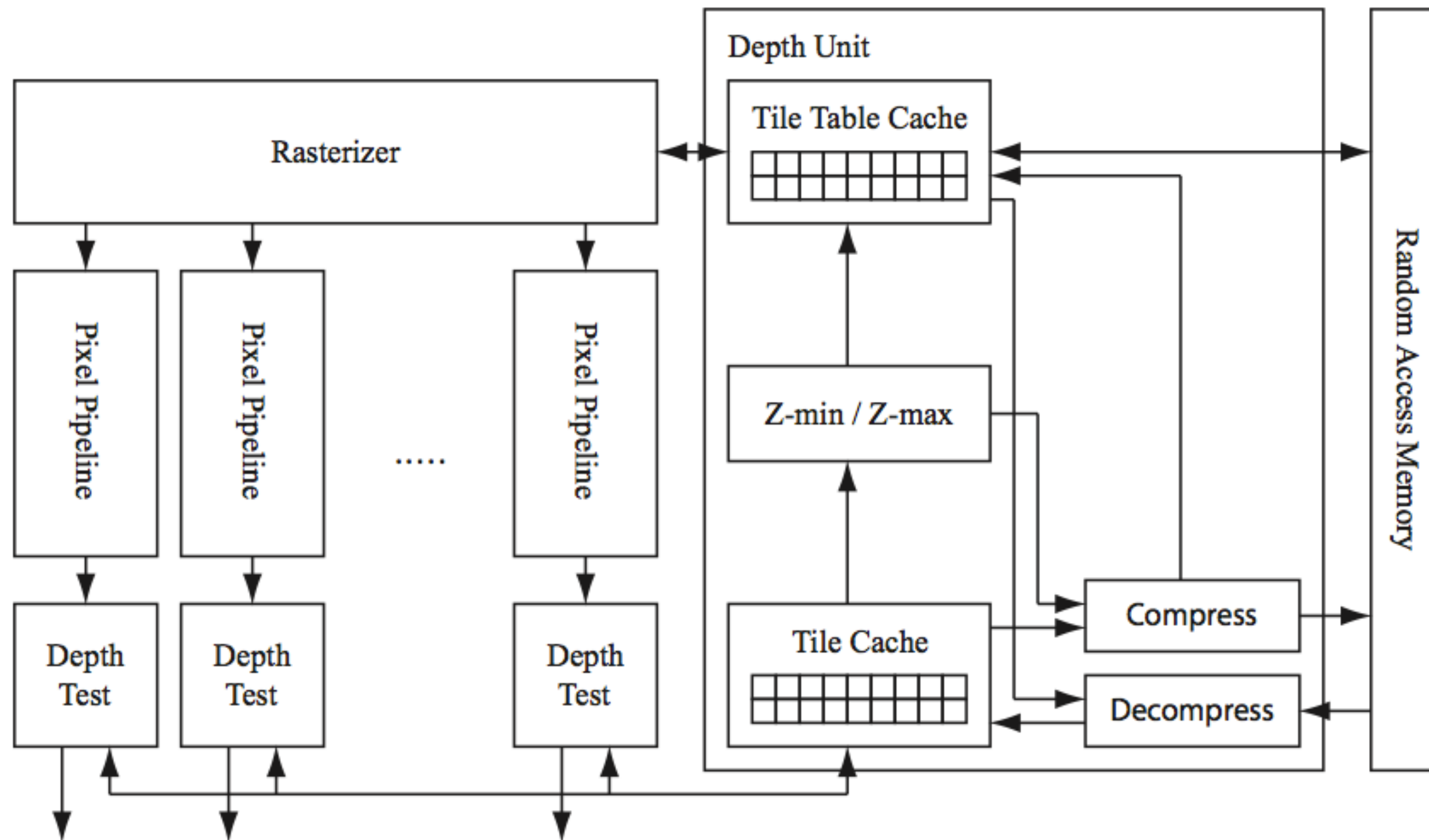
Frame-buffer compression

Depth-buffer compression

- **Motivation: reduce bandwidth required for depth-buffer accesses**
 - **Worst-case (uncompressed) buffer allocated in DRAM**
 - **Conserving memory footprint is a non-goal**
(Need for real-time guarantees in graphics applications requires application to plan for worst case anyway)
- **Lossless compression**
 - **Q. Why not lossy?**
- **Designed for fixed-point numbers (fixed-point math in rasterizer)**

Depth-buffer compression is tile based

- Main idea: exploit similarity of values within a screen tile



On tile evict:

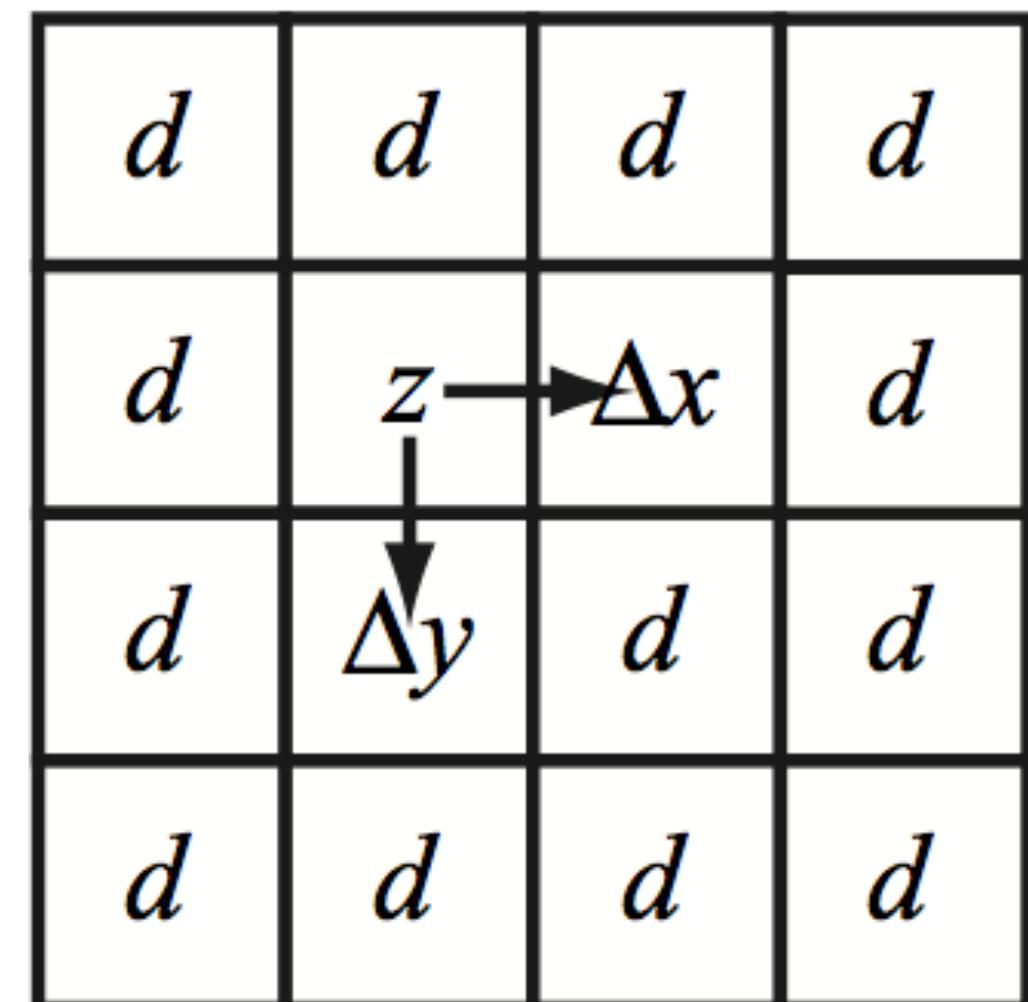
1. Compute zmin/zmax (needed for hierarchical culling and/or compression)
2. Attempt to compress
3. Update tile table
4. Store tile to memory

On tile load:

1. Check tile table for compression scheme
2. Load required bits from memory
3. Decompress into tile cache

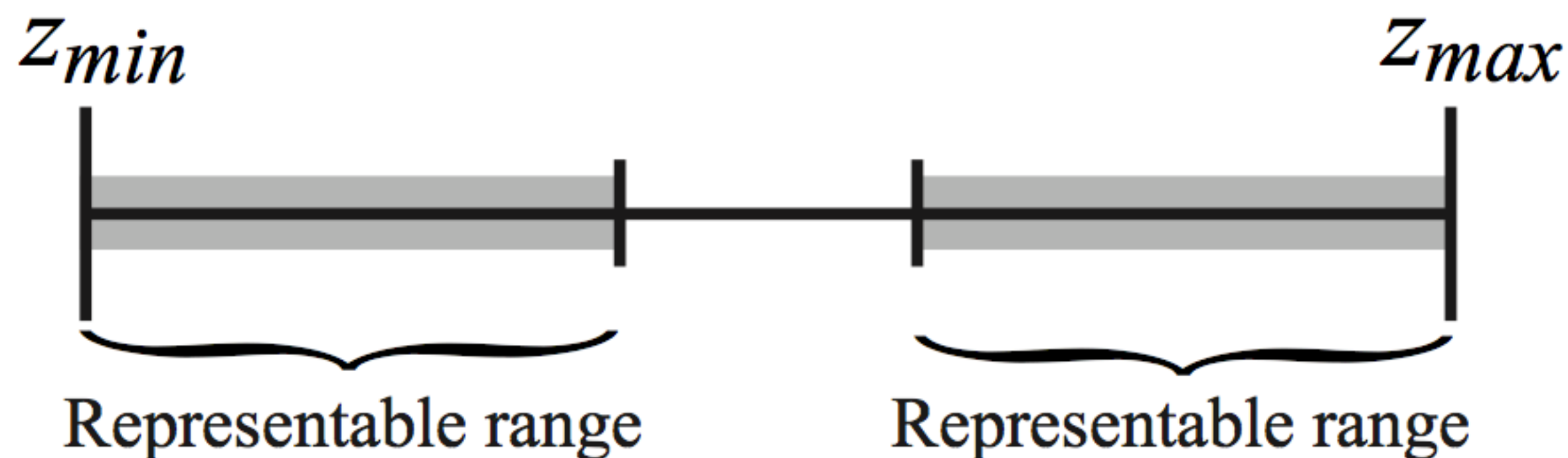
Anchor encoding

- Choose anchor value and compute DX, DY from adjacent pixels (fits a plane to the data)
- Use plane to predict depths at other pixels, store offset d from prediction at each pixel
- Scheme (for 24-bit depth buffer)
 - Anchor: 24 bits (full resolution)
 - DX, DY: 15 bits
 - Per-sample offsets: 5 bits



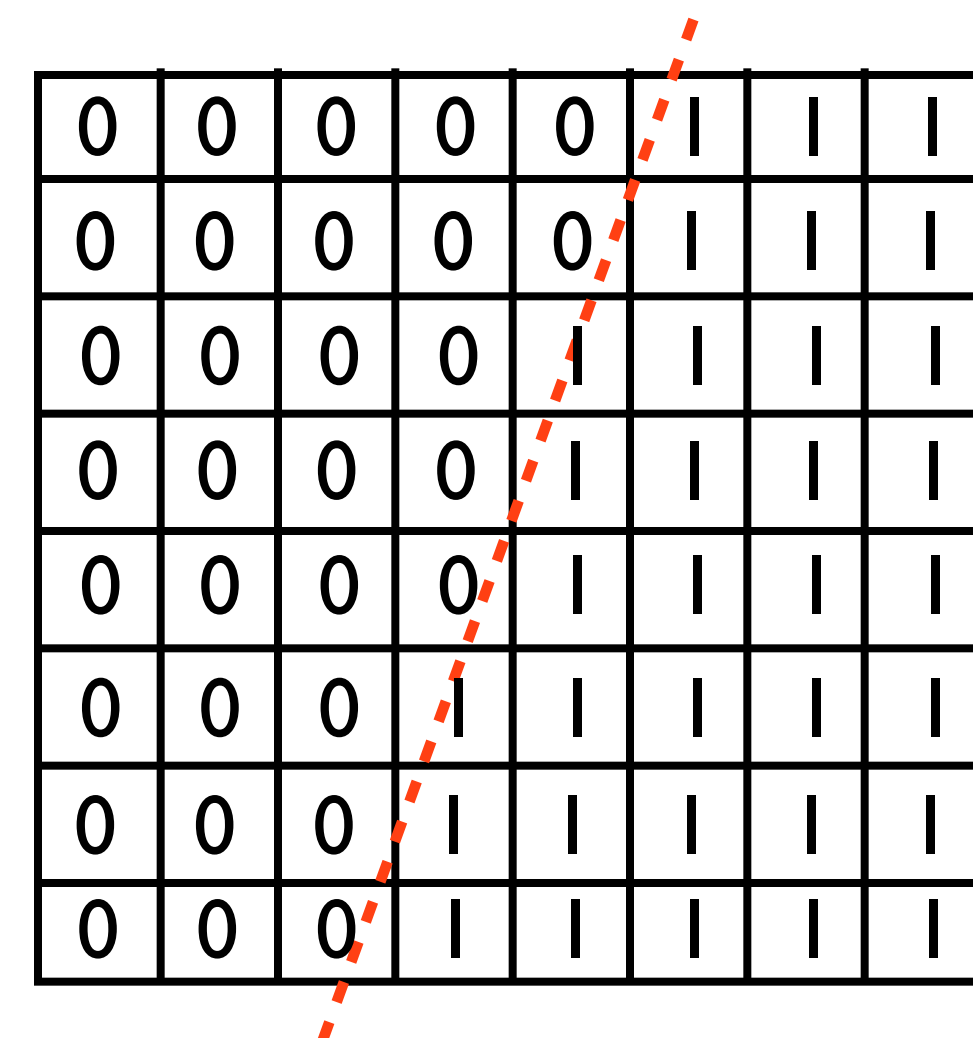
Depth-offset compression

- Assume depth values have low dynamic range relative to tile's z_{min} and z_{max} (assume two surfaces)
- Store z_{min}/z_{max} (need to anyway for hierarchical Z)
- Store low-precision (8-12 bits) offset value for each sample
 - MSB encodes if offset is from z_{min} or z_{max}



Explicit plane encoding

- **Do not attempt to infer prediction plane, just get the plane equation directly from the rasterizer**
 - Store plane equation in tile (values must be stored with high precision: to match exact math performed by rasterizer)
 - Store bit per sample indicating coverage
- **Simple extension to multiple triangles per tile:**
 - Store up to N plane equations in tile
 - Store $\log_2(N)$ bit id per depth sample indicating which triangle it belongs to
- **When new triangle contributes coverage to tile:**
 - Add new plane equation if storage is available, else decompress
- **To decompress:**
 - For each sample, evaluate $Z(x,y)$ for appropriate plane



0	0	0	0	0	1	1	1
0	0	0	0	0	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	1	1	1	1	1
0	0	0	1	1	1	1	1
0	0	0	1	1	1	1	1

Summary: reducing the bandwidth requirements of depth testing

- **Caching: access DRAM less often (by caching depth buffer data)**
 - **Hierarchical Z techniques (zmin/zmax culling): “early outs” result in accesses individual sample data less often**
 - **Data compression: reduce number of bits that must be read from memory**
-
- **Color buffer is also compressed using similar techniques**
 - **Depth buffer typically achieves higher compression ratios than color buffer. Why?**

Visibility/culling relationships

- **Hierarchical traversal during rasterization**
 - Leveraged to reduce coverage testing and occlusion work
 - Tile size likely coupled to hierarchical Z granularity
 - May also be coupled to compression tile granularity
- **Hierarchical culling and plane-based buffer compression are most effective when triangles are reasonably large (recall triangle size discussion in lecture 2)**

Stochastic rasterization

■ Accurate camera simulation in real-time rendering

- Visibility algorithms discussed today simulate image formation by a virtual pinhole camera with and infinitely fast shutter
- Real cameras have finite apertures and finite exposure duration
- Accurate camera simulation requires visibility computation to be integrated over time and lens aperture (very high computational cost, see readings)

Time integration: motion blur



Lens aperture integration: defocus blur



Readings

- **M. Abrash, Rasterization on Larrabee, Dr. Dobbs Portal. 5/1/2009**