# Computation in Big Spaces

John Skilling (john@skilling.co.uk)      AI in Astronomy 2019

The 20th century got a lot wrong.

1. It argued about Bayes.
2. It viewed functions through 19th century lens.
3. It viewed inference through Laplace transform.
4. It thought dimensionality was hard.

Result:   Principles of inference were obscured.
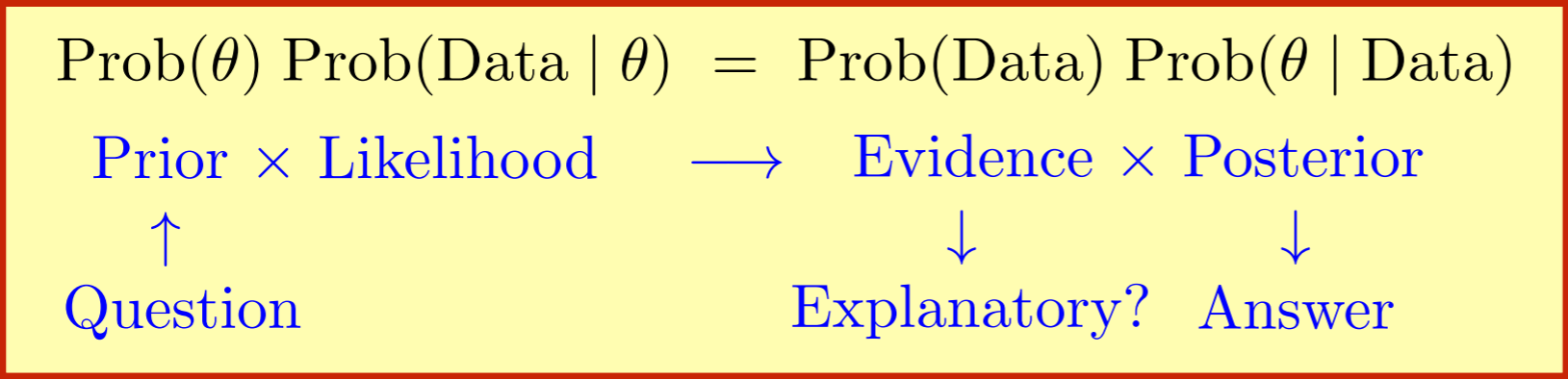          Easy stuff was judged impossible.

**It's 2019: we can do better!**

# Bayes

Fact: the only consistent calculus for uncertainty is ordinary probability.
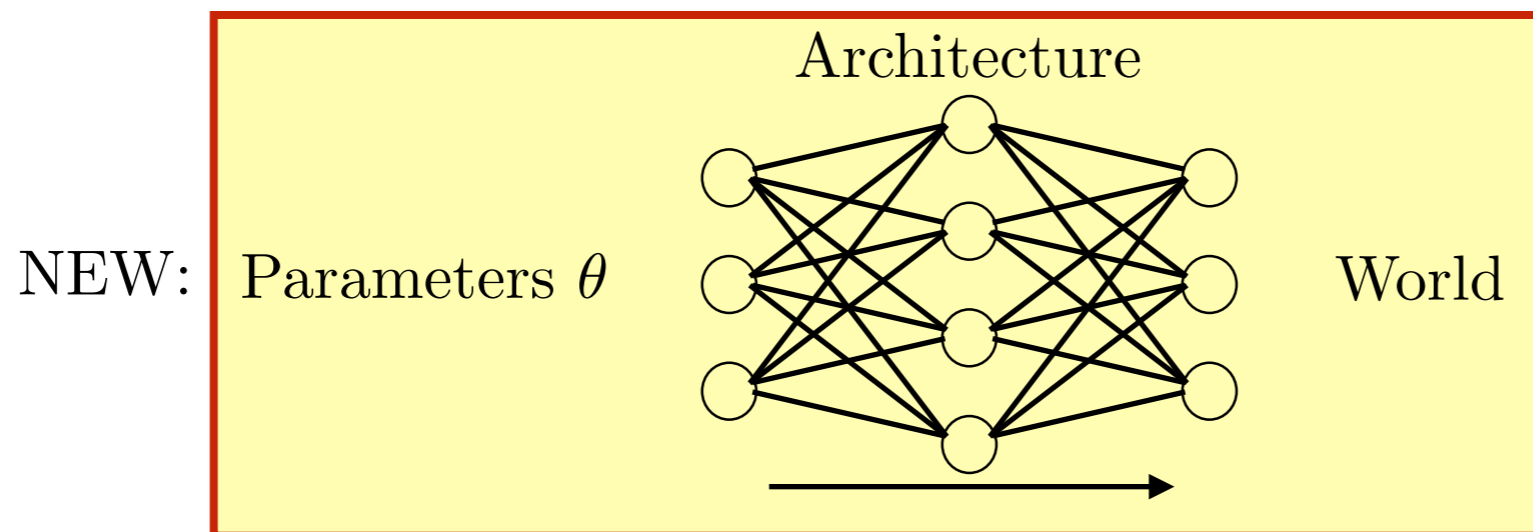
Sum rule:        proportions (probabilities) add up.     ✓✓

Product rule:    proportions (probabilities) multiply.

$$\text{Prob}(\theta)\,\text{Prob}(\text{Data}\,|\,\theta) \;=\; \text{Prob}(\text{Data})\,\text{Prob}(\theta\,|\,\text{Data})$$

Prior $\times$ Likelihood    $\longrightarrow$    Evidence $\times$ Posterior

$\uparrow$                        $\downarrow$         $\downarrow$

Question              Explanatory?   Answer

**THE** framework
of inference

Prior: *What's your world model?*     (human crafted, deep-learned, …)

OLD:   World $= \sum \int \left( \begin{array}{c} \text{modified Hankel functions} \\ \text{of 3rd type} \ldots\ldots \text{etc.} \end{array} \right)$ (parameters $\theta$)

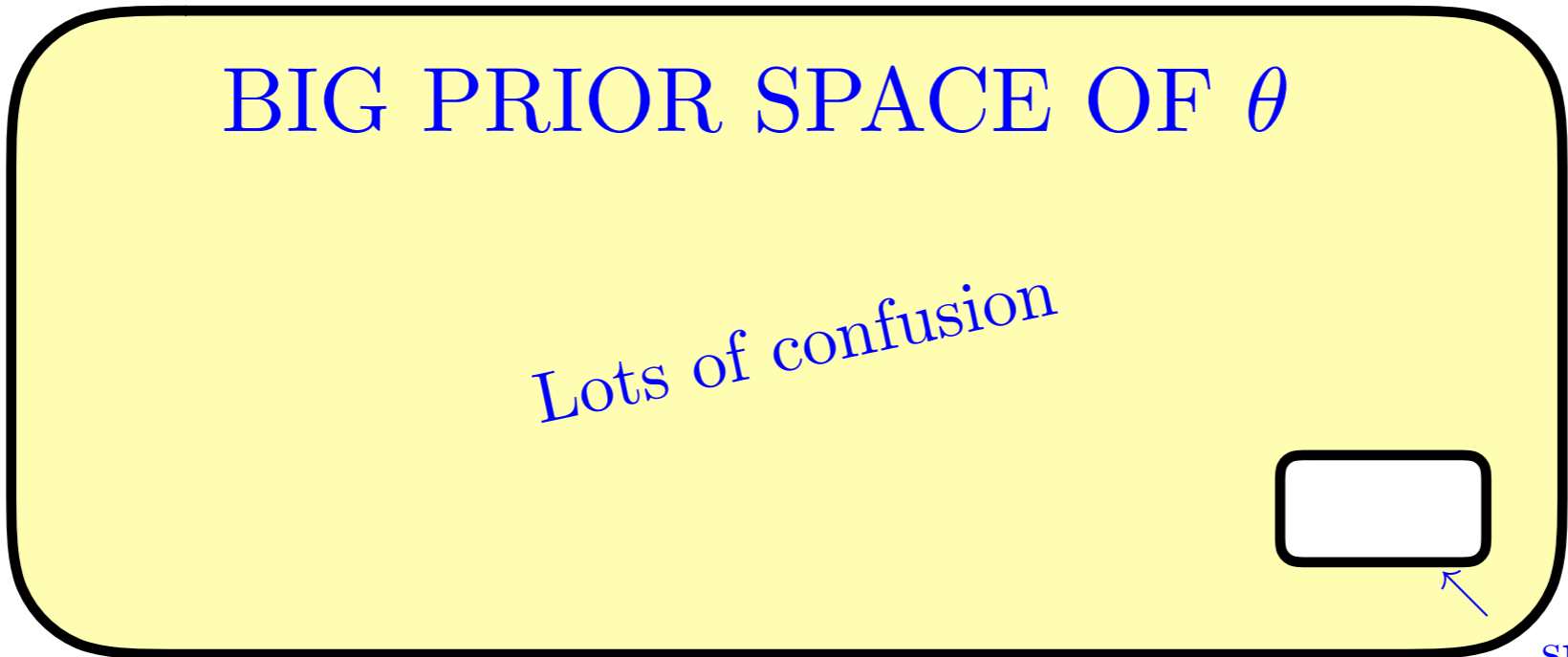Architecture

NEW:   Parameters $\theta$                   World     ✓✓

# Computational Inference

$$\text{Prob}(\theta)\,\text{Prob}(\text{Data}\,|\,\theta) \;=\; \text{Prob}(\text{Data})\,\text{Prob}(\theta\,|\,\text{Data})$$

$$d\pi(\theta) \times L(\theta) \qquad \longrightarrow \qquad Z \;\times\; dP(\theta)$$

$$\uparrow \qquad\qquad\qquad\qquad\qquad \downarrow \qquad\quad \downarrow$$

Question                          Explanatory?   Answer

**THE** framework
of inference

(1)  Evidence:  $Z = \int L(\theta)\, d\pi(\theta)$     ⟵ Looks hard but isn't.

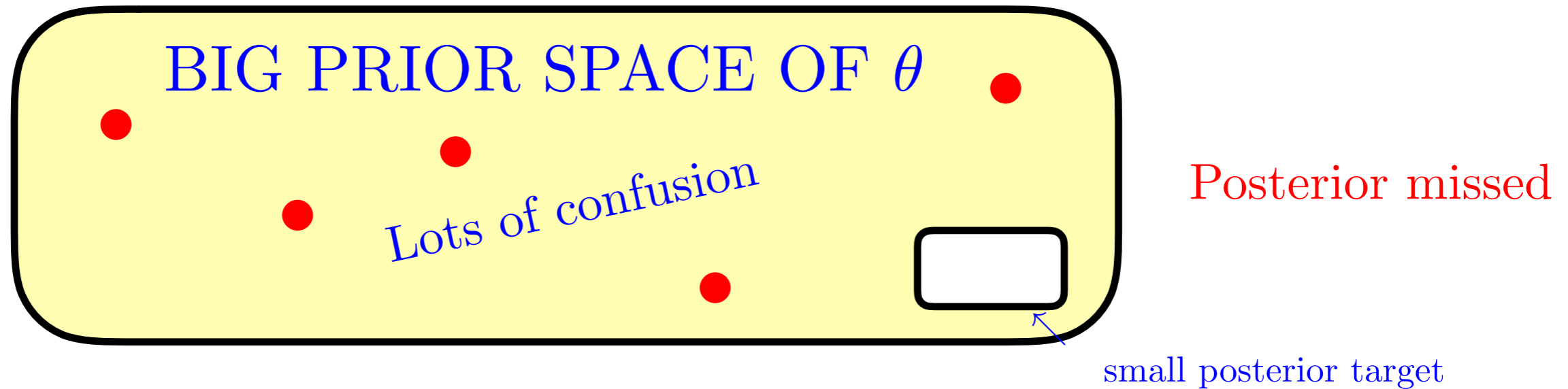(2)  Posterior:  $dP(\theta) = \dfrac{L(\theta)\, d\pi(\theta)}{Z}$     ⟵ Secondary.

## BIG PRIOR SPACE OF $\theta$

Lots of confusion

small posterior target

Can't look everywhere. Must sample (Monte Carlo).

# Hopeless methods for $\quad Z = \int L(\theta) \, d\pi(\theta)$

Direct: $\qquad\qquad Z = \left\langle L \right\rangle_{\text{prior}}$

BIG PRIOR SPACE OF $\theta$

Lots of confusion

Posterior missed

small posterior target

# Hopeless methods for $\quad Z = \int L(\theta)\, d\pi(\theta)$

Direct: $\qquad\qquad Z = \big\langle L \big\rangle_{\text{prior}}$



BIG PRIOR SPACE OF $\theta$

Lots of confusion

Posterior missed

small posterior target

Harmonic mean: $\quad \dfrac{1}{Z} = \Big\langle \dfrac{1}{L} \Big\rangle_{\text{posterior}}$



BIG PRIOR SPACE OF $\theta$

Lots of confusion

Coverage fraction $\Delta\pi$ lost

small posterior target

Must connect prior-to-posterior.

Must connect prior-to-posterior.  ✓ ✓

---

**Another bad idea (annealing)**

$$Z(\beta) = \int L(\theta)^\beta \, d\pi(\theta)$$

$\beta = 0$ prior $\longrightarrow$ $\beta = 1$ posterior

---

The 20th century got a lot wrong.

1. It argued about Bayes.
2. It viewed functions through 19th century lens.
3. It viewed inference through Laplace transform.
4. It thought dimensionality was hard.
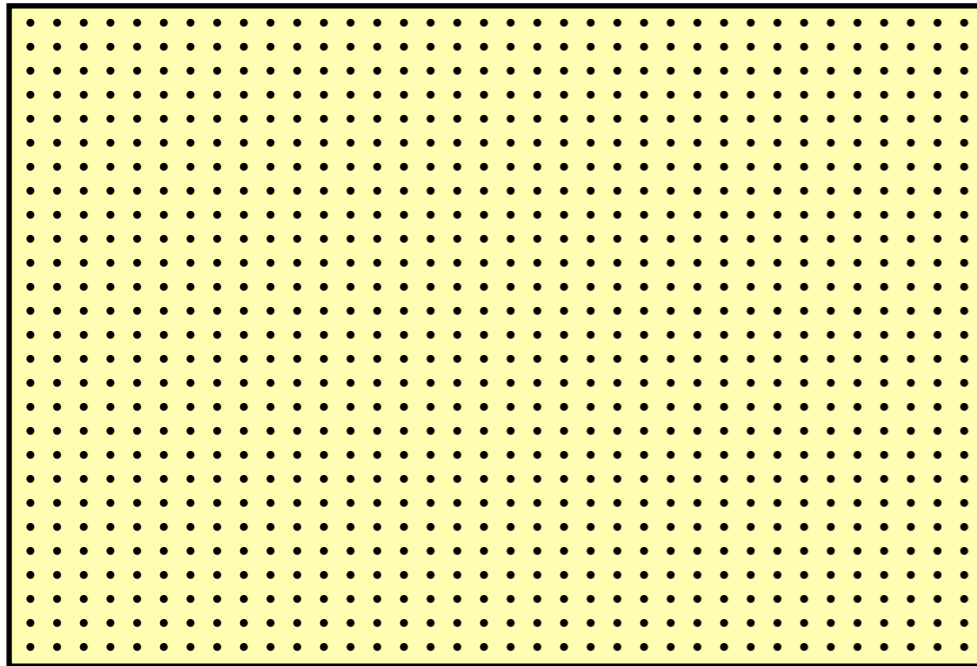
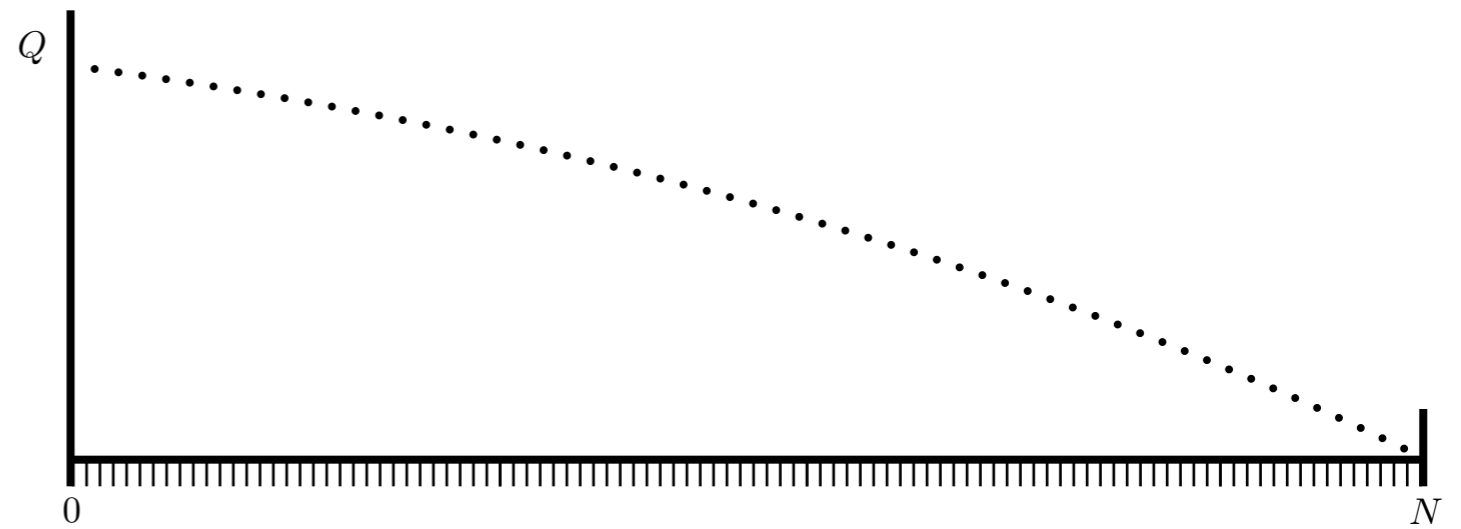**Back to first principles — SMASH DIMENSIONALITY!**

✓ ✓

Q:    How does a mathematician find a needle in a haystack?

A:    Keep halving the haystack and discarding the "wrong" half.
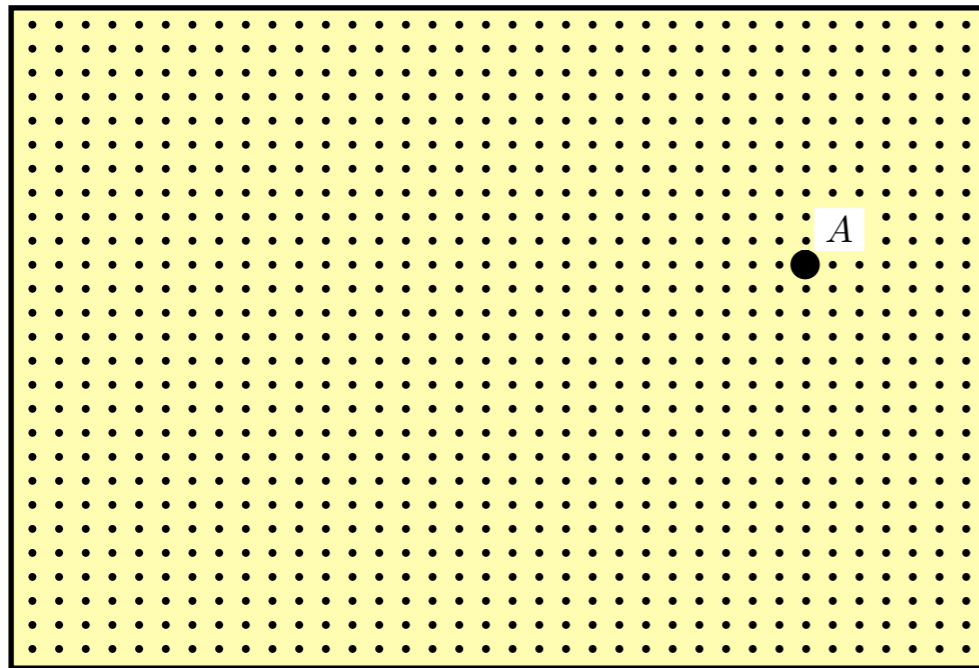      Performance: factor of 2 per step.



$N$ possible points

$N$ points ranked by quality

7

## Compression

Q: How does a programmer find a small target?
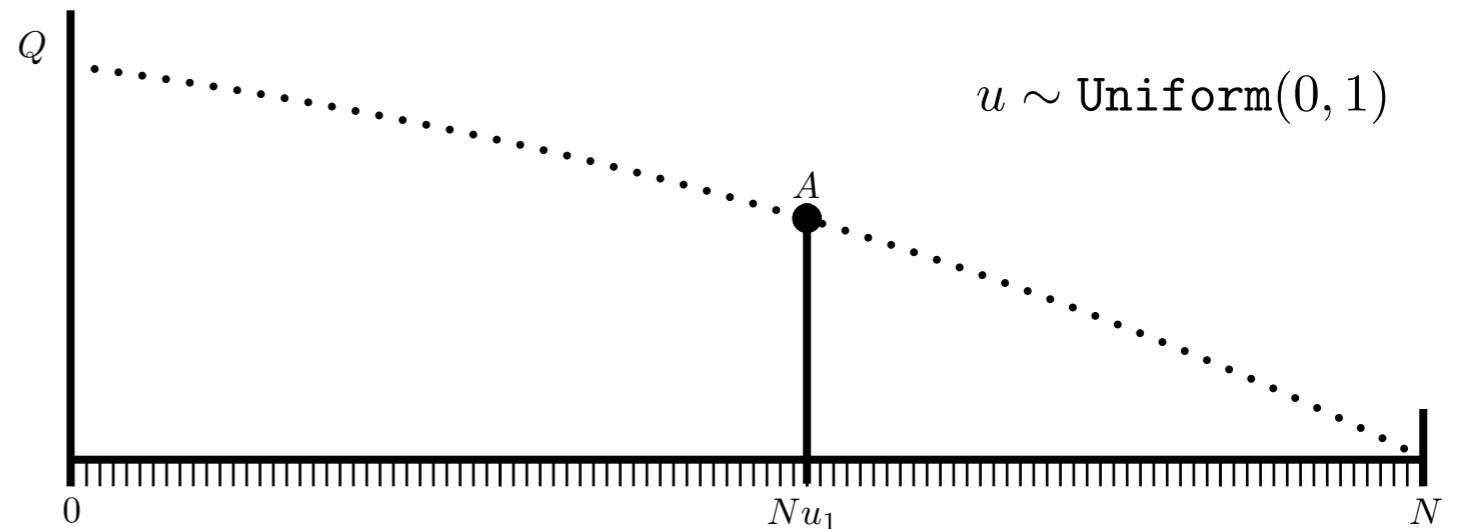A: Keep taking a random point and discarding everywhere worse.
Performance: statistically a factor of $e$ per step.



$A$ random

$\log(u_1) = -1 \pm 1$

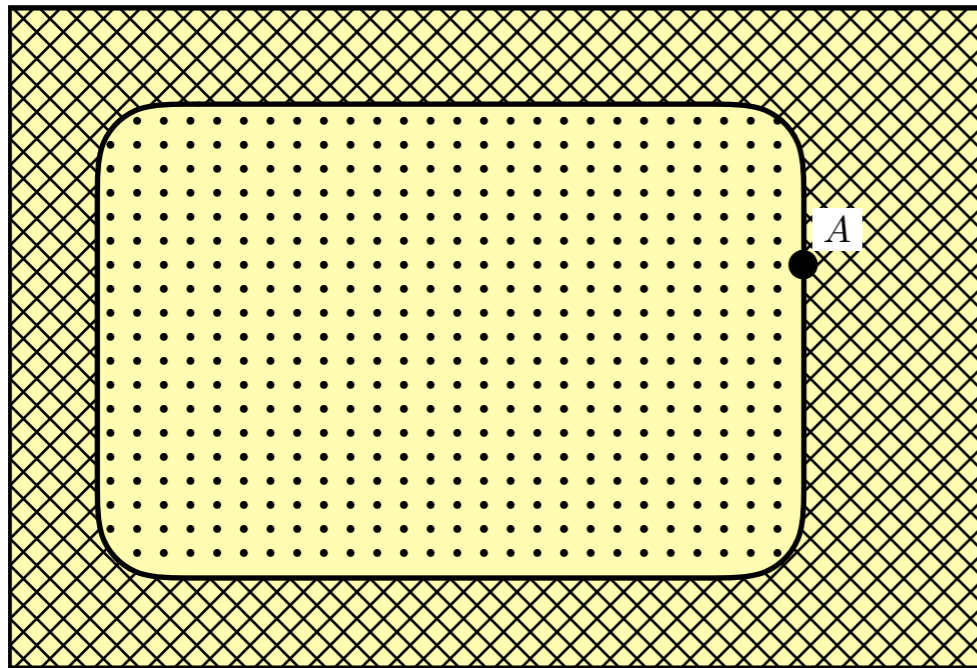$u \sim \texttt{Uniform}(0,1)$

$N$ possible points

$N$ points ranked by quality

## Compression

Q:  How does a programmer find a small target?

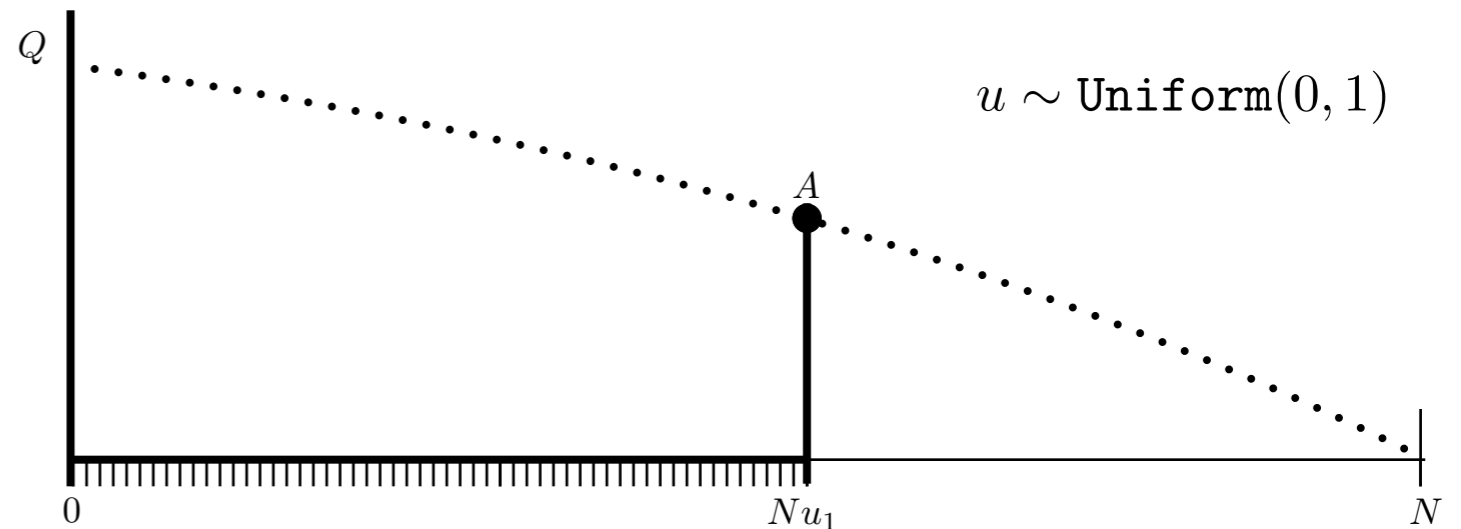A:  Keep taking a random point and discarding everywhere worse.

Performance: statistically a factor of $e$ per step.



$A$ random

$\log(u_1) = -1 \pm 1$

$u \sim \texttt{Uniform}(0,1)$

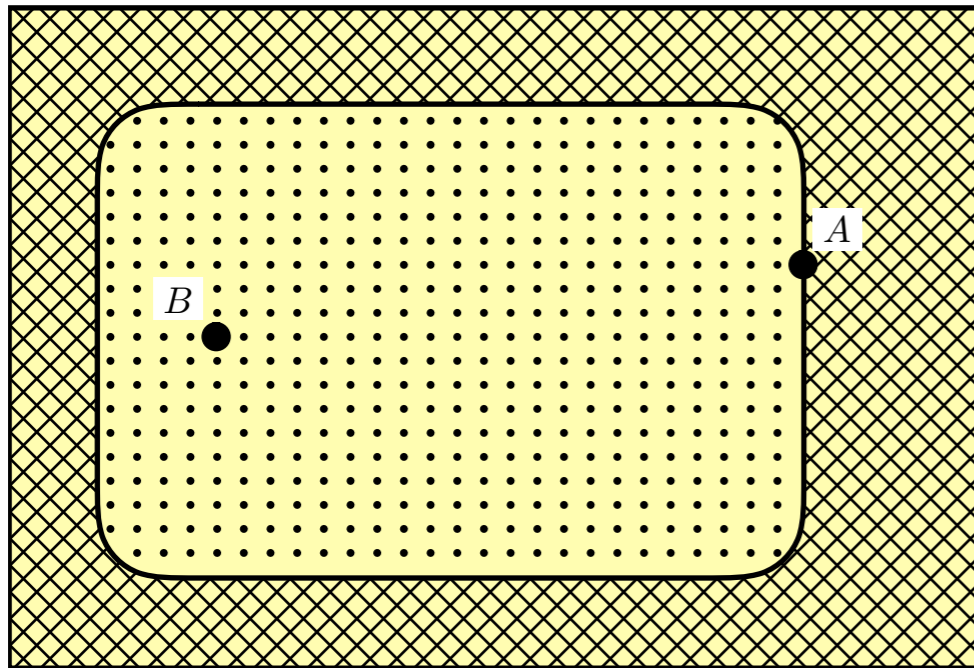$N$ possible points

$N$ points ranked by quality

## Compression
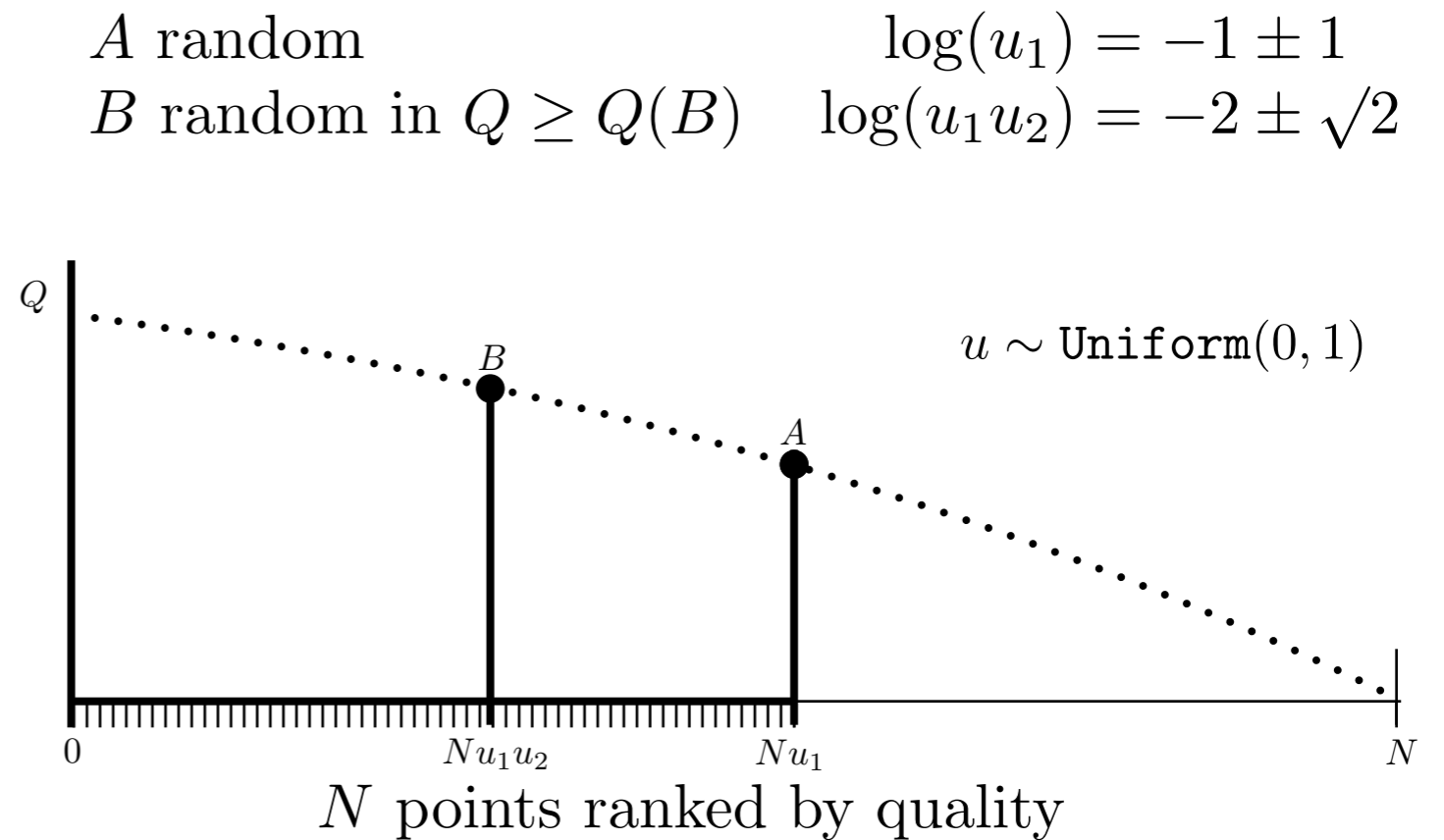
Q:    How does a programmer find a small target?

A:    Keep taking a random point and discarding everywhere worse.
      Performance: statistically a factor of $e$ per step.

$A$ random                           $\log(u_1) = -1 \pm 1$

$B$ random in $Q \geq Q(B)$          $\log(u_1 u_2) = -2 \pm \sqrt{2}$



$u \sim \texttt{Uniform}(0,1)$
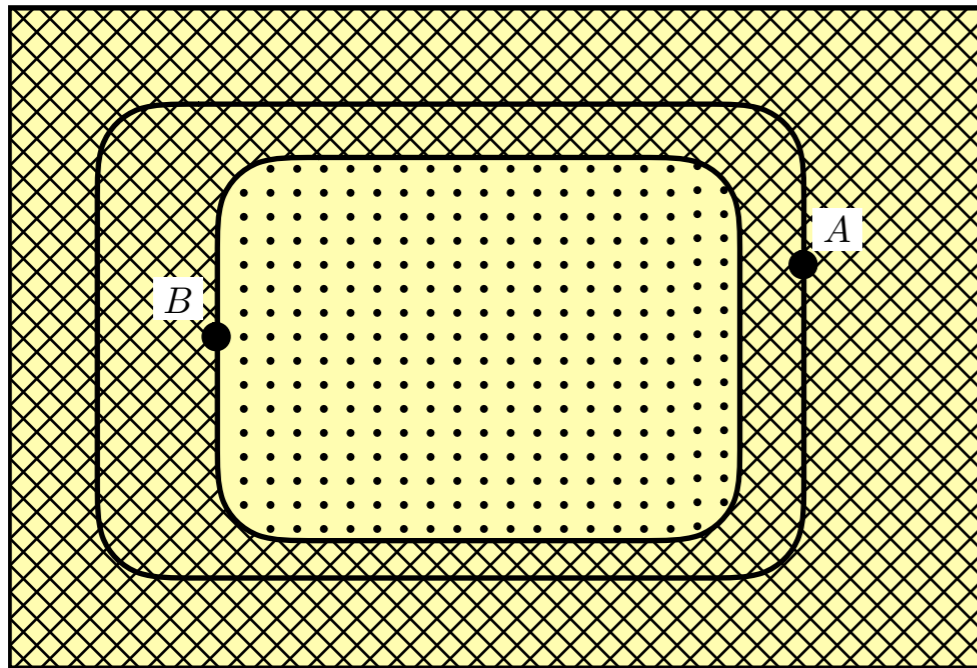
$N$ possible points
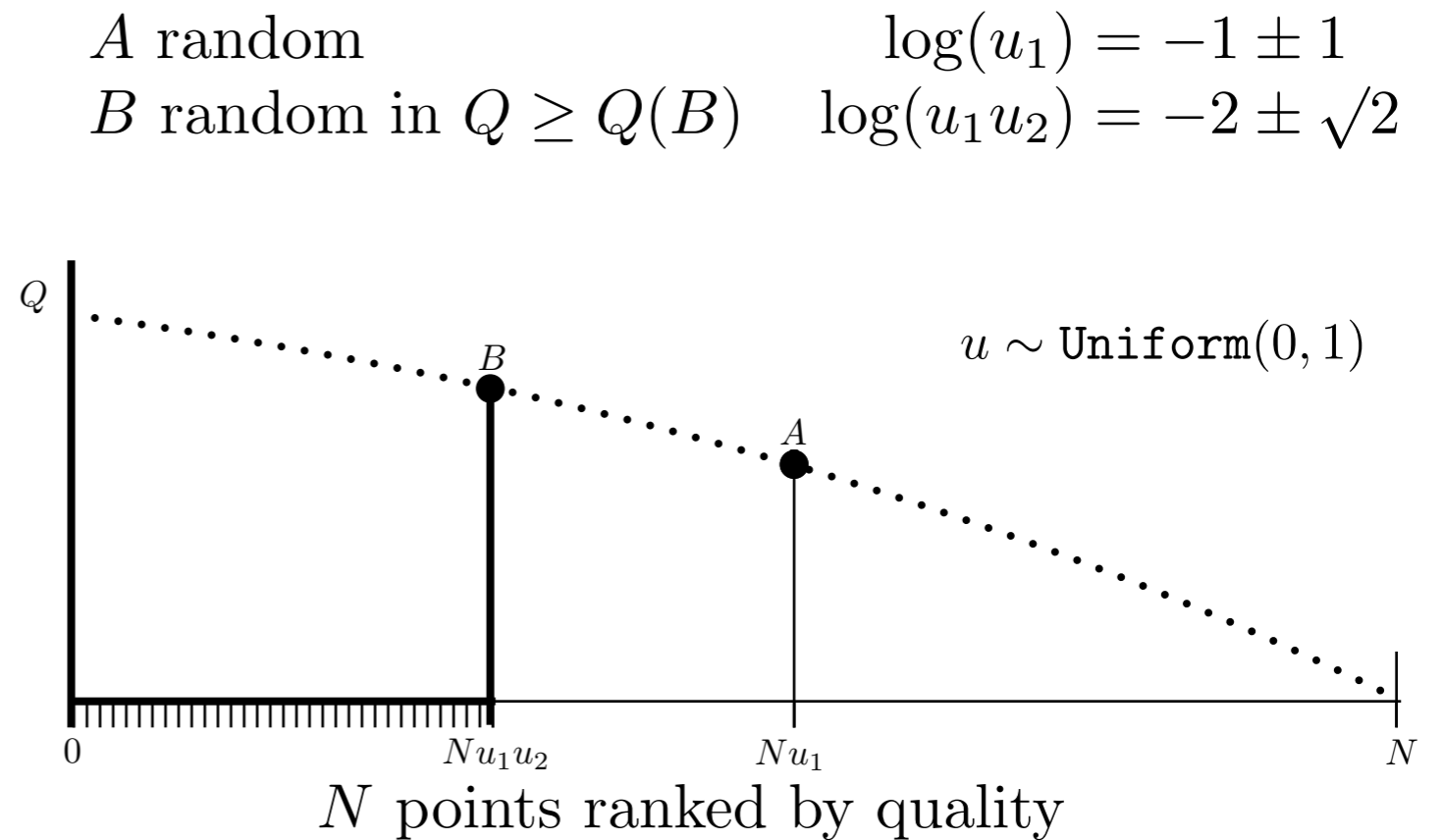
$N$ points ranked by quality

10

Q:   How does a programmer find a small target?
A:   Keep taking a random point and discarding everywhere worse.
     Performance: statistically a factor of $e$ per step.



$A$ random                          $\log(u_1) = -1 \pm 1$
$B$ random in $Q \geq Q(B)$         $\log(u_1 u_2) = -2 \pm \sqrt{2}$

$u \sim \texttt{Uniform}(0,1)$

$N$ possible points
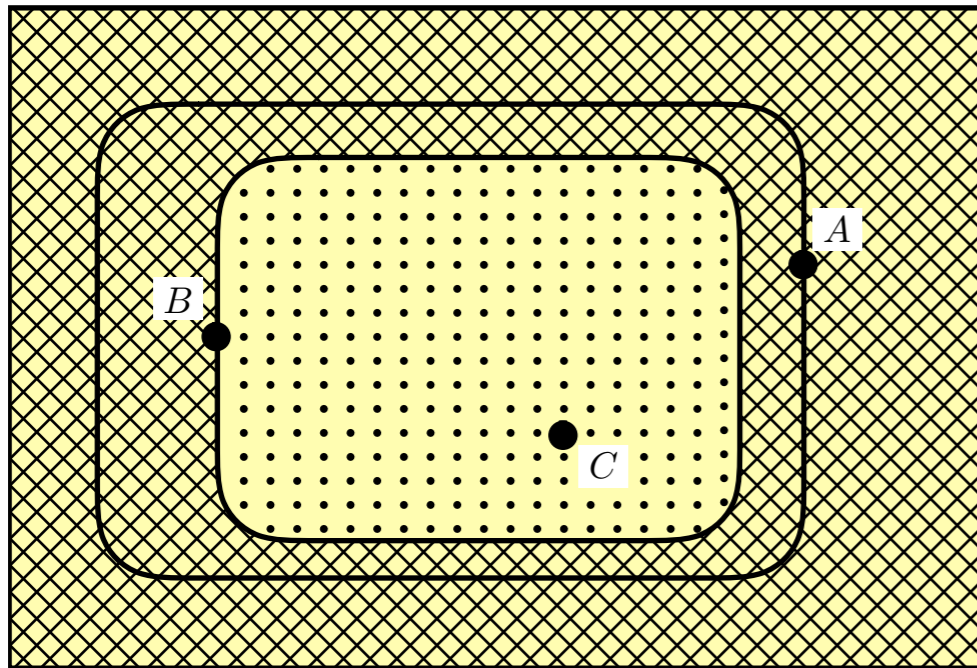
$N$ points ranked by quality

## Compression
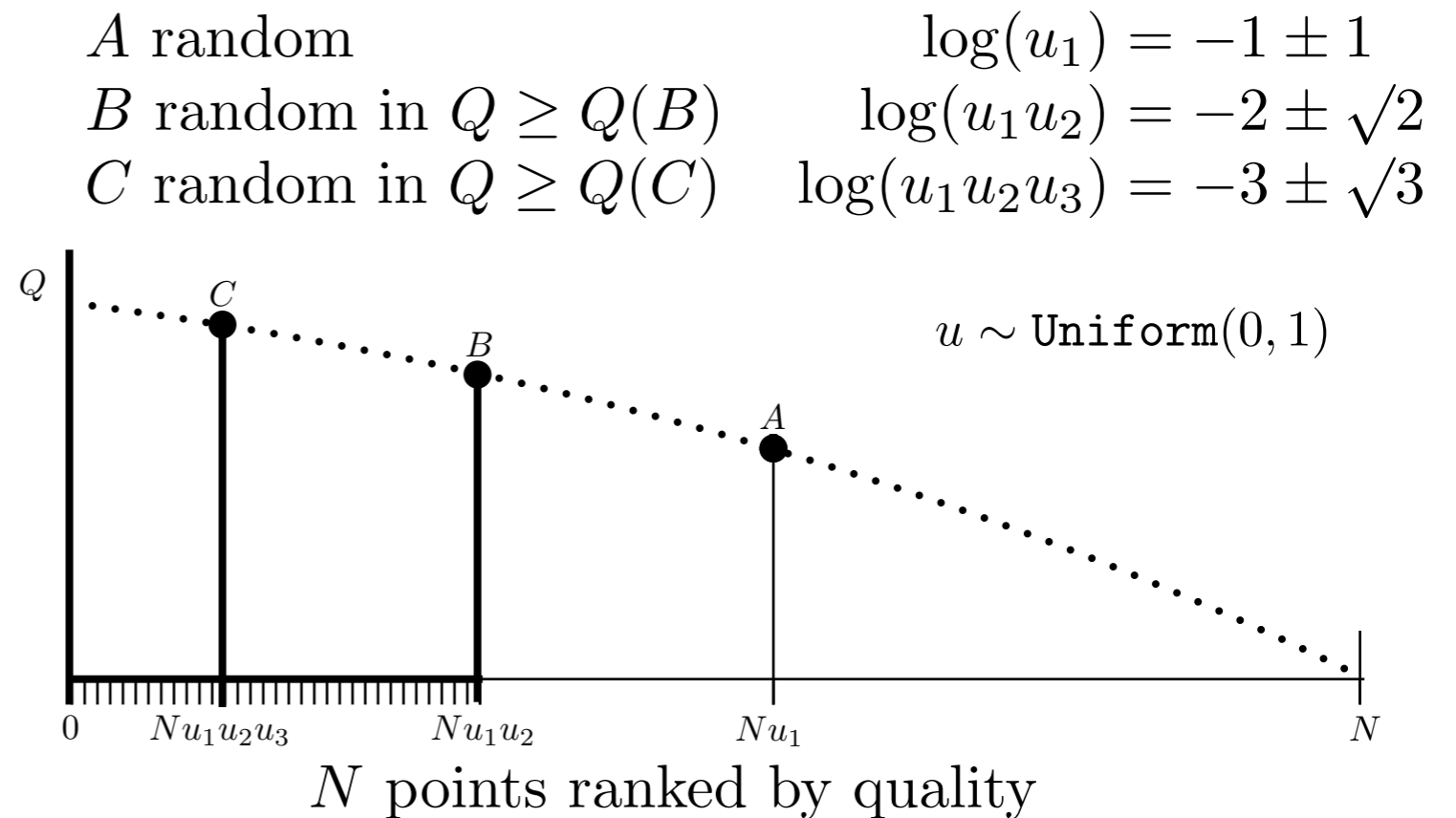
Q:    How does a programmer find a small target?

A:    Keep taking a random point and discarding everywhere worse.
       Performance: statistically a factor of $e$ per step.



$N$ possible points

$A$ random                    $\log(u_1) = -1 \pm 1$
$B$ random in $Q \geq Q(B)$    $\log(u_1 u_2) = -2 \pm \sqrt{2}$
$C$ random in $Q \geq Q(C)$    $\log(u_1 u_2 u_3) = -3 \pm \sqrt{3}$

$u \sim \texttt{Uniform}(0, 1)$

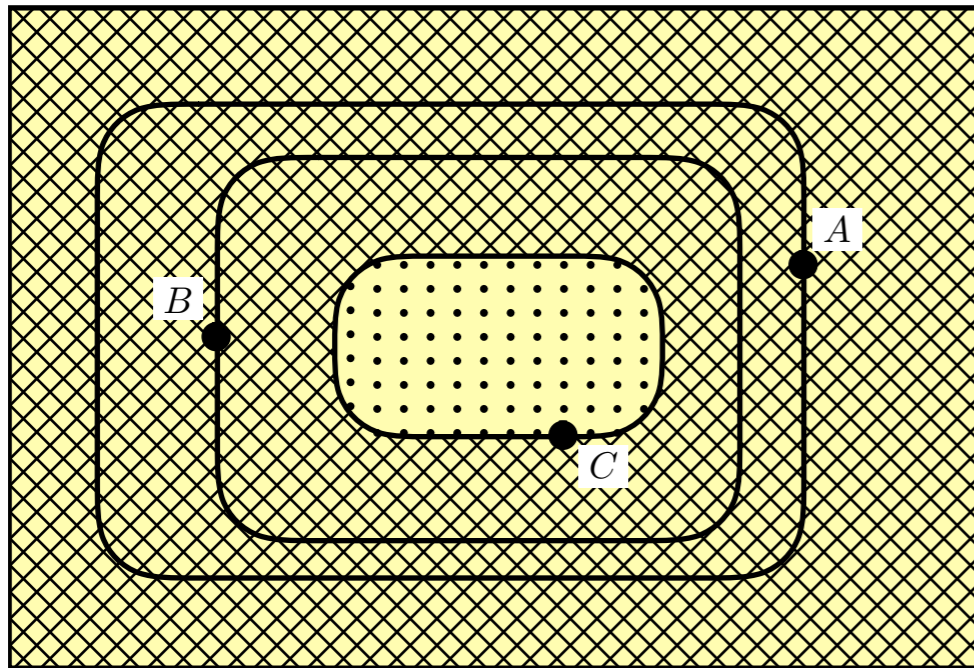$N$ points ranked by quality

# Compression

Q:  How does a programmer find a small target?

A:  Keep taking a random point and discarding everywhere worse.

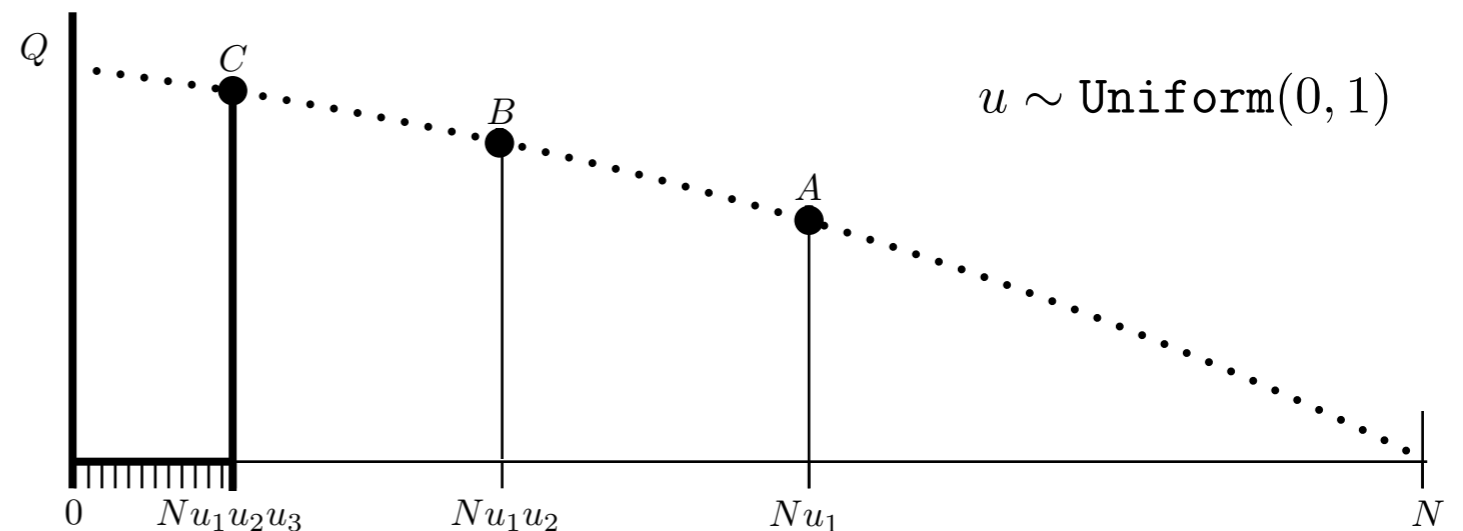Performance: statistically a factor of $e$ per step.

$A$ random

$B$ random in $Q \geq Q(B)$

$C$ random in $Q \geq Q(C)$

$\log(u_1) = -1 \pm 1$

$\log(u_1 u_2) = -2 \pm \sqrt{2}$

$\log(u_1 u_2 u_3) = -3 \pm \sqrt{3}$

$u \sim \texttt{Uniform}(0, 1)$

$N$ possible points

$N$ points ranked by quality

13

## Compression

Q: How does a programmer find a small target?
A: Keep taking a random point and discarding everywhere worse.
   Performance: statistically a factor of $e$ per step.



$A$ random

$B$ random in $Q \geq Q(B)$

$C$ random in $Q \geq Q(C)$

$\log(u_1) = -1 \pm 1$

$\log(u_1 u_2) = -2 \pm \sqrt{2}$

$\log(u_1 u_2 u_3) = -3 \pm \sqrt{3}$
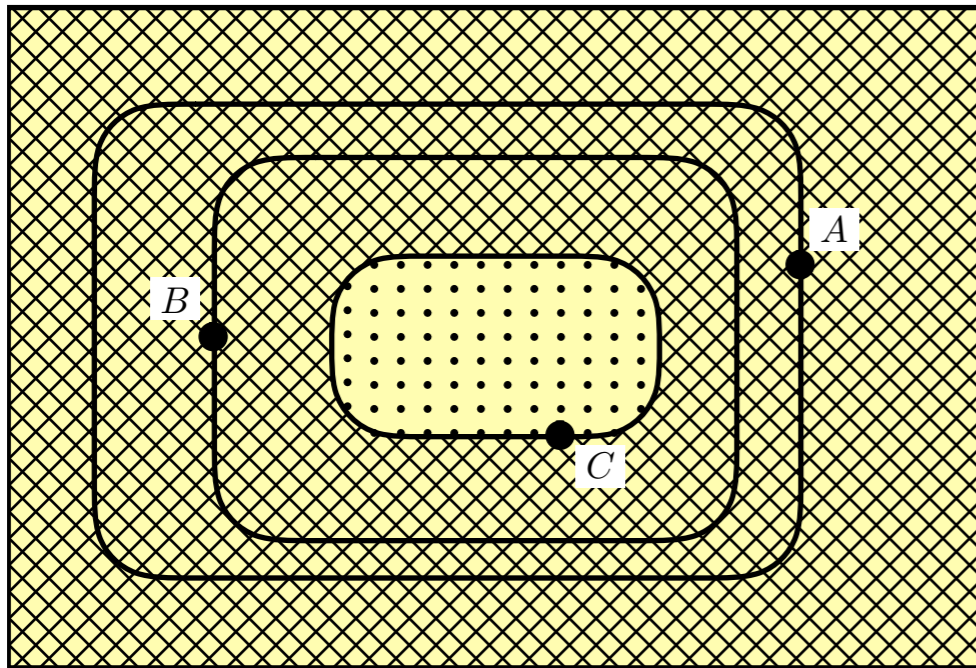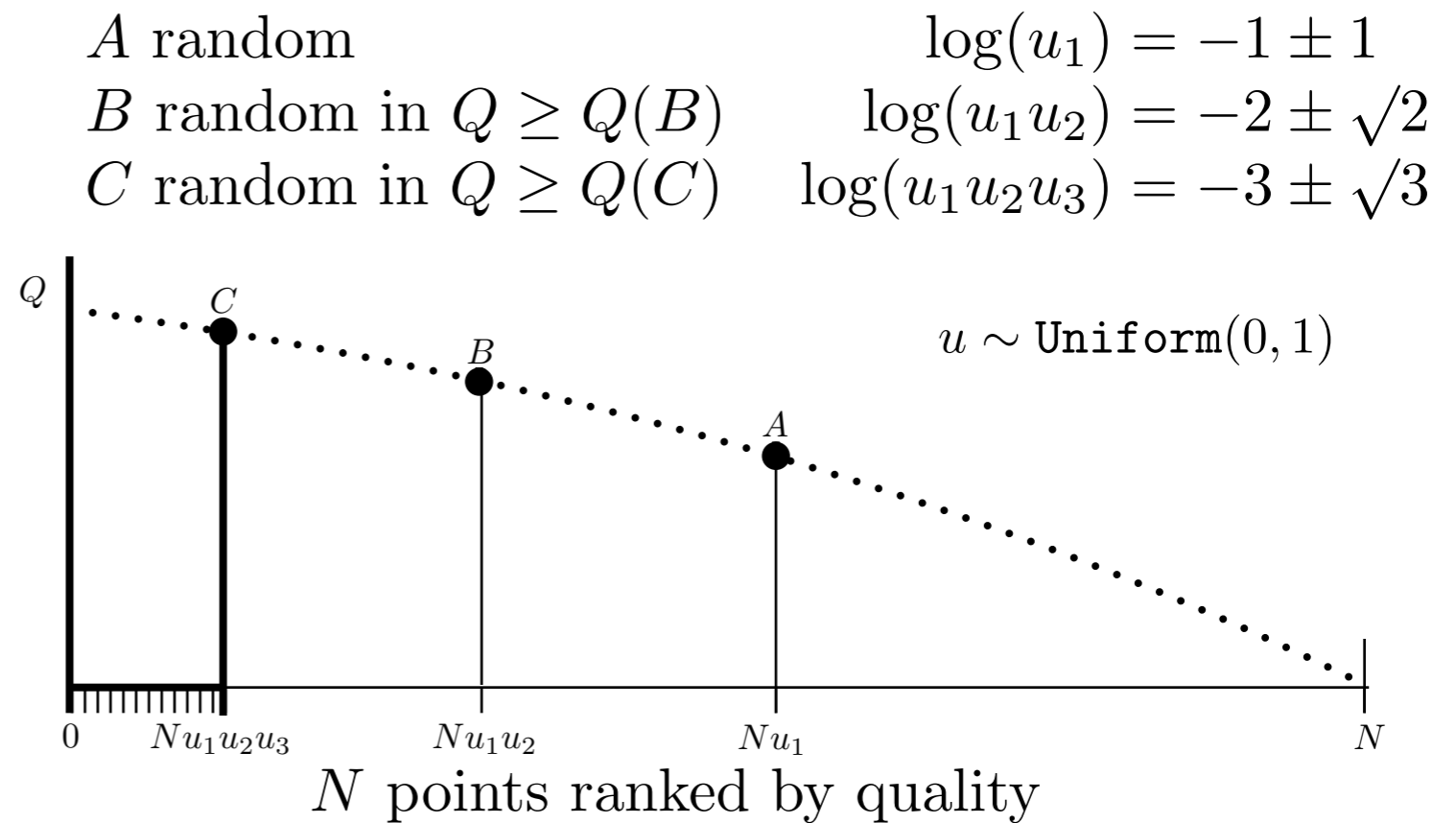
$u \sim \mathtt{Uniform}(0,1)$

$N$ possible points

$N$ points ranked by quality

After $k$ steps, accumulated compression $e^{k \pm \sqrt{k}}$ has enclosed quality $Q \geq Q_k$.

$$\log\left(\frac{\#\text{ targets}}{\#\text{ possibles}}\right) = -k \pm \sqrt{k}$$

Compression estimated statistically in log(ratio) steps *without exploring everywhere*.
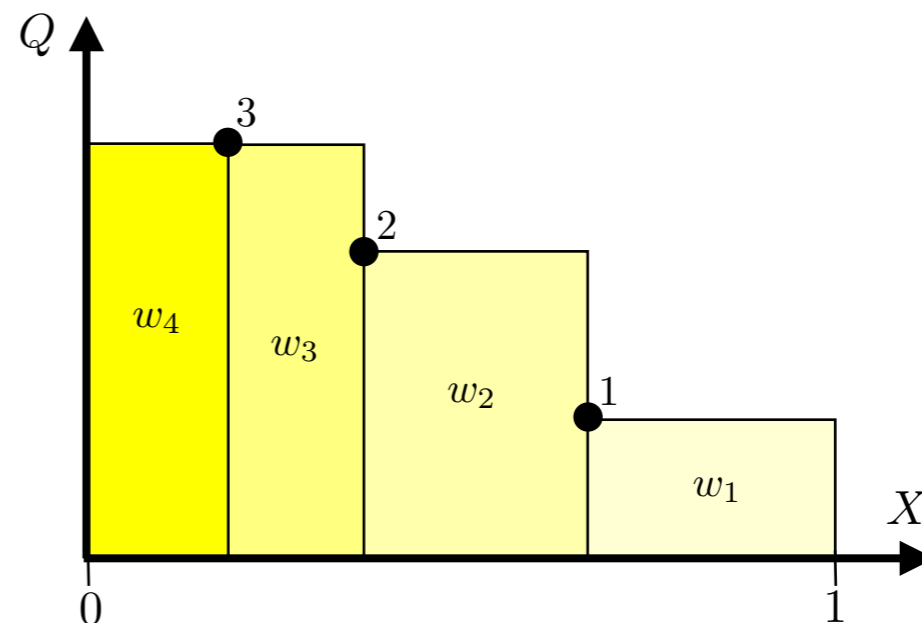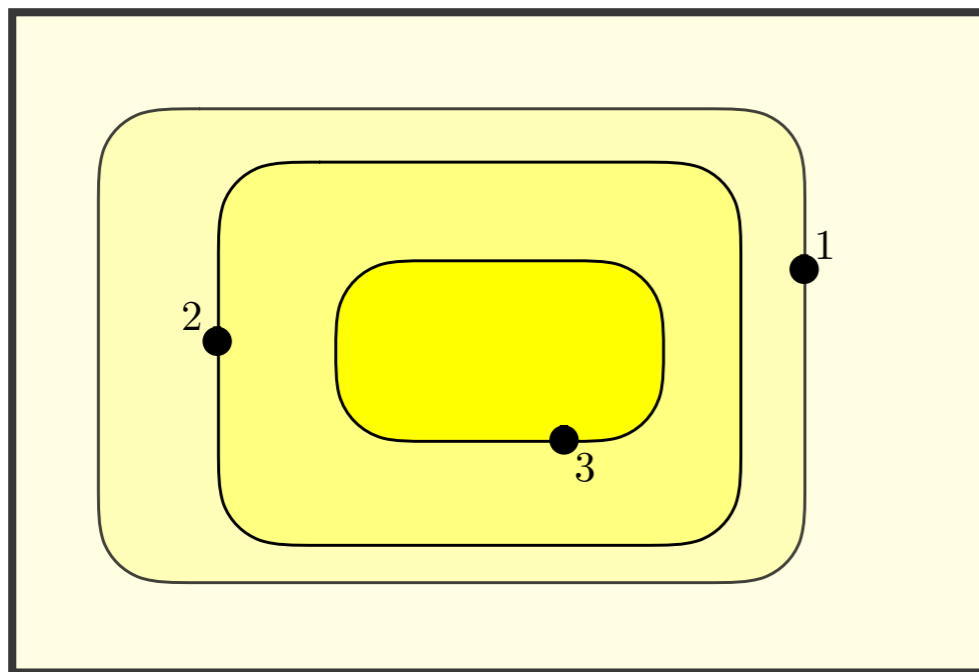
This is the *nested sampling* algorithm.

## Quantification

A nested sampling run yields a sequence $k = 1, 2, 3, \ldots$ of

$$\begin{cases} \mathbf{x}_k = \text{location} \\ Q(\mathbf{x}_k) = \text{quality} \\ X_k = \text{fraction of possibilities with quality} \geq Q, \text{ as statistical estimate} \end{cases}$$

Thus it gives the relationship $Q(X)$ and a sample location $\mathbf{x}$ in each shell of quality.



Shell $k$ contributes $\quad w_k = Q_k \, \Delta X_k \quad$ to $\quad \displaystyle\int Q(\mathbf{x})d\mathbf{x} = \int Q(X)dX \approx \sum Q_k \, \Delta X_k.$
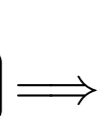
So we can integrate in any geometrical dimension — we just need the background measure.

## Quantification — Bayes

Quality = Likelihood

$$\underbrace{\text{Prob}(\mathbf{x})}_{\text{Prior } dX} \times \underbrace{\text{Prob}(\text{data} \mid \mathbf{x})}_{\text{Likelihood } L} = \text{Prob}(\mathbf{x}, \text{data}) = \underbrace{\text{Prob}(\text{data})}_{\text{Evidence } Z} \times \underbrace{\text{Prob}(\mathbf{x} \mid \text{data})}_{\text{Posterior } dP}$$
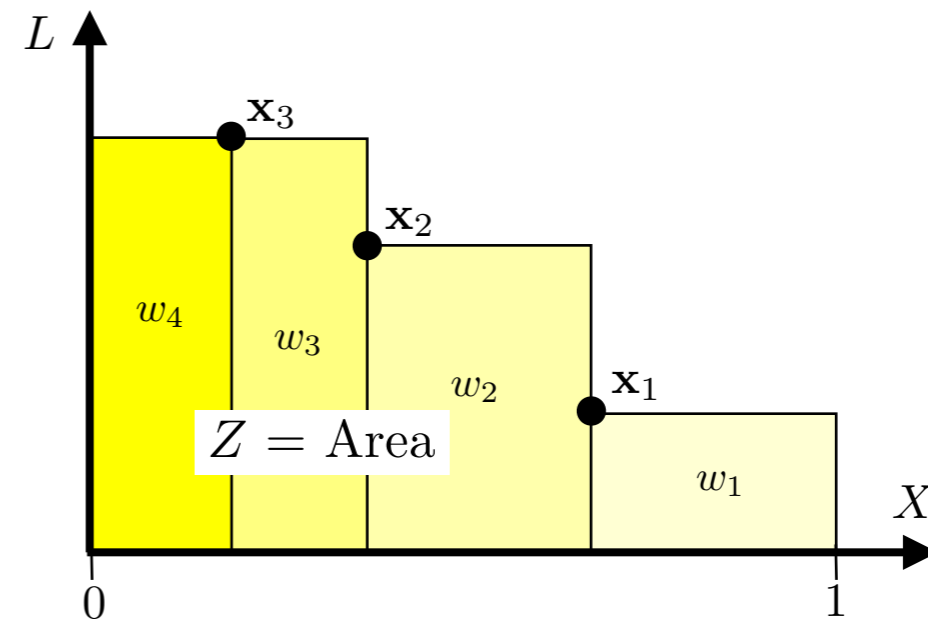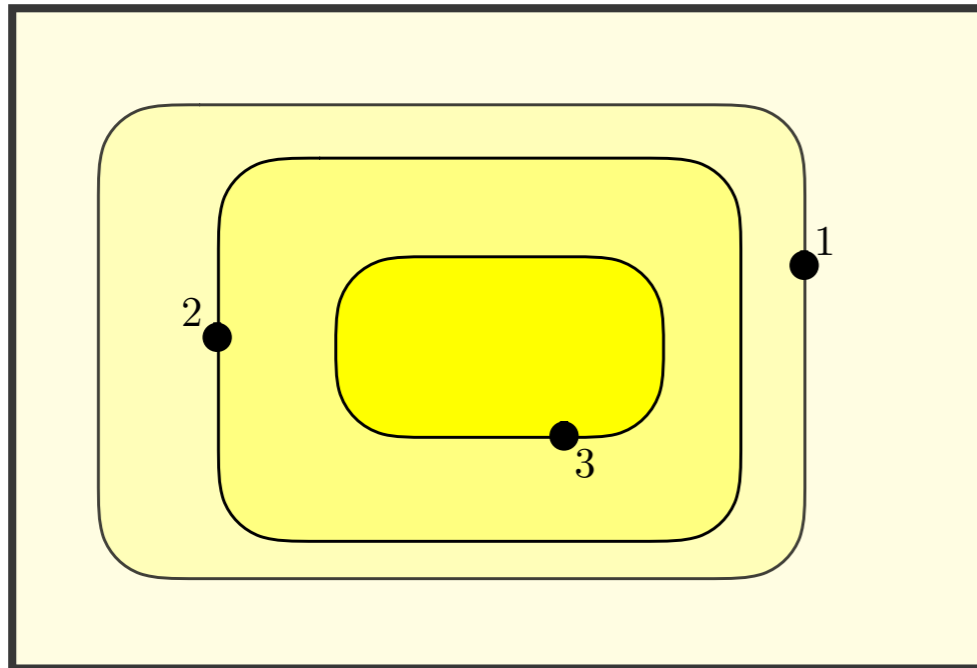
Inputs $\Longrightarrow$

$\begin{cases} \text{Prior} & \text{is} \quad \pi(\mathbf{x}), \quad \text{equivalent to uniform on } 0 < X < 1 \\[2mm] \text{Likelihood} & \text{is} \quad L(\mathbf{x}) \quad \text{or} \quad L(X) \\[2mm] \text{Evidence} & \text{is} \quad Z = \int L(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} = \int_0^1 L(X)dX \quad \text{or} \quad Z \approx \sum w_k \qquad (1) \\[2mm] \text{Posterior} & \text{is} \quad \text{weighted samples } \mathbf{x}_1, \mathbf{x}_2, \ldots, \text{ with} \quad \text{Prob}(k) = w_k/Z \qquad (2) \end{cases}$ $\Longrightarrow$ Outputs

Compression yields posterior samples on the fly.
Unified computation of evidence and posterior.

# Take-home message:

Bayes is required for consistent inference and
it's not as hard as you may have thought.

Not all that glitters is gold:

| | |
|---|---|
| Maximum . . . | *non-Bayesian* |
| Posterior | *semi-Bayesian* |
| Posterior + Evidence | *Bayesian* |

John Skilling, AI for Astronomy 2019

Garching-bei-München, EU