

Extreme Memoization Everything in a LUT!

Pratyush Patel
Luis Ceze

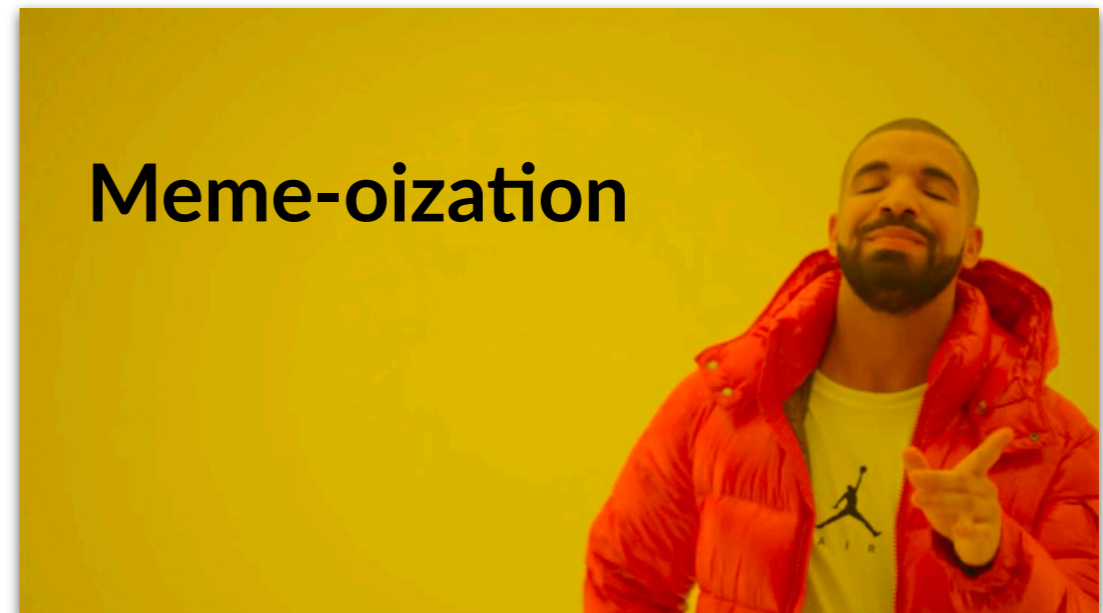


Extreme *Meme*-oization Everything in a LUT!

Pratyush Patel
Luis Ceze



Extreme *Meme*-oization Everything in a LUT!





DENNARD
SCALING
1974—2006

MOORE'S
LAW
1965—??

RIP

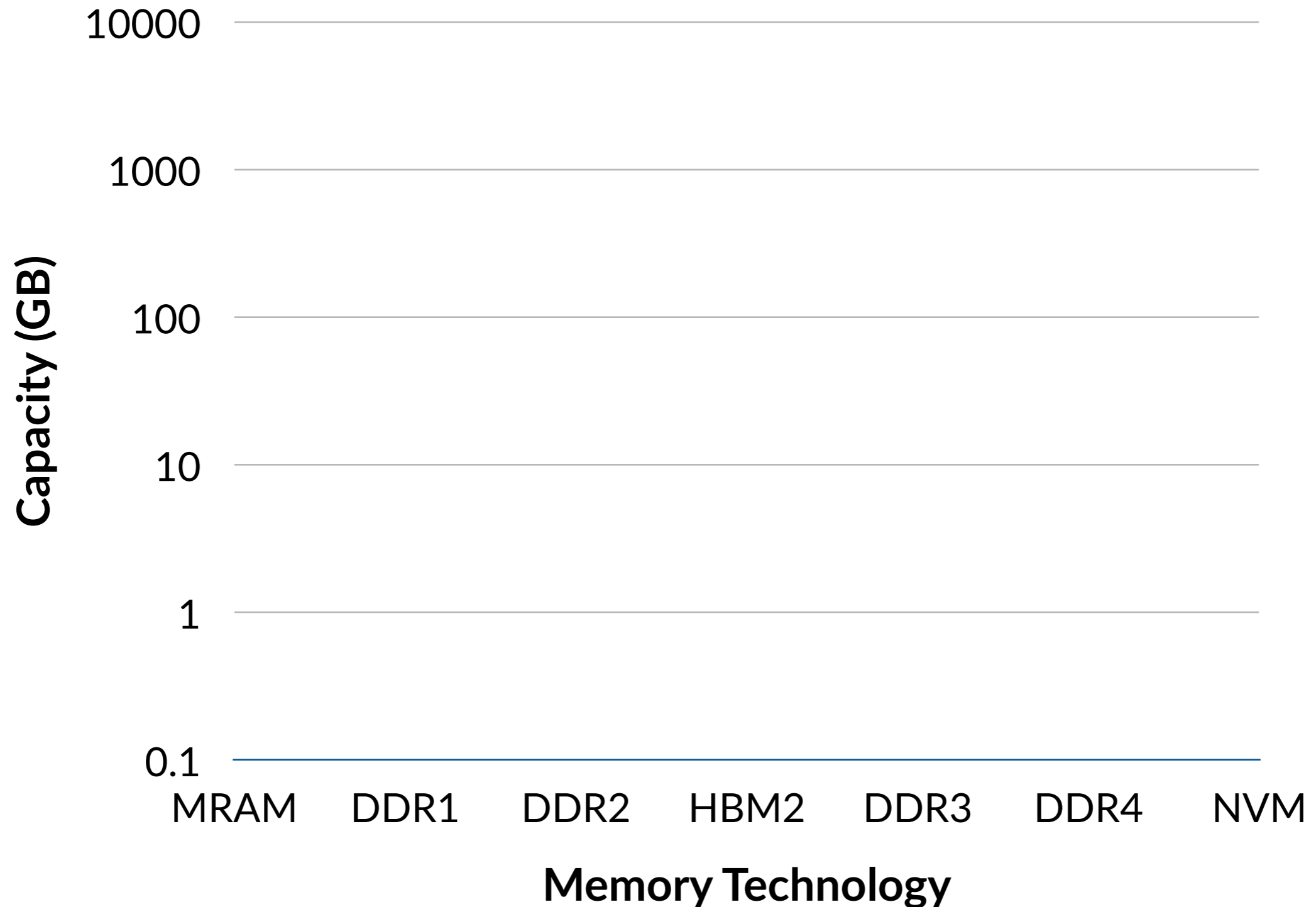
RIP

RIP

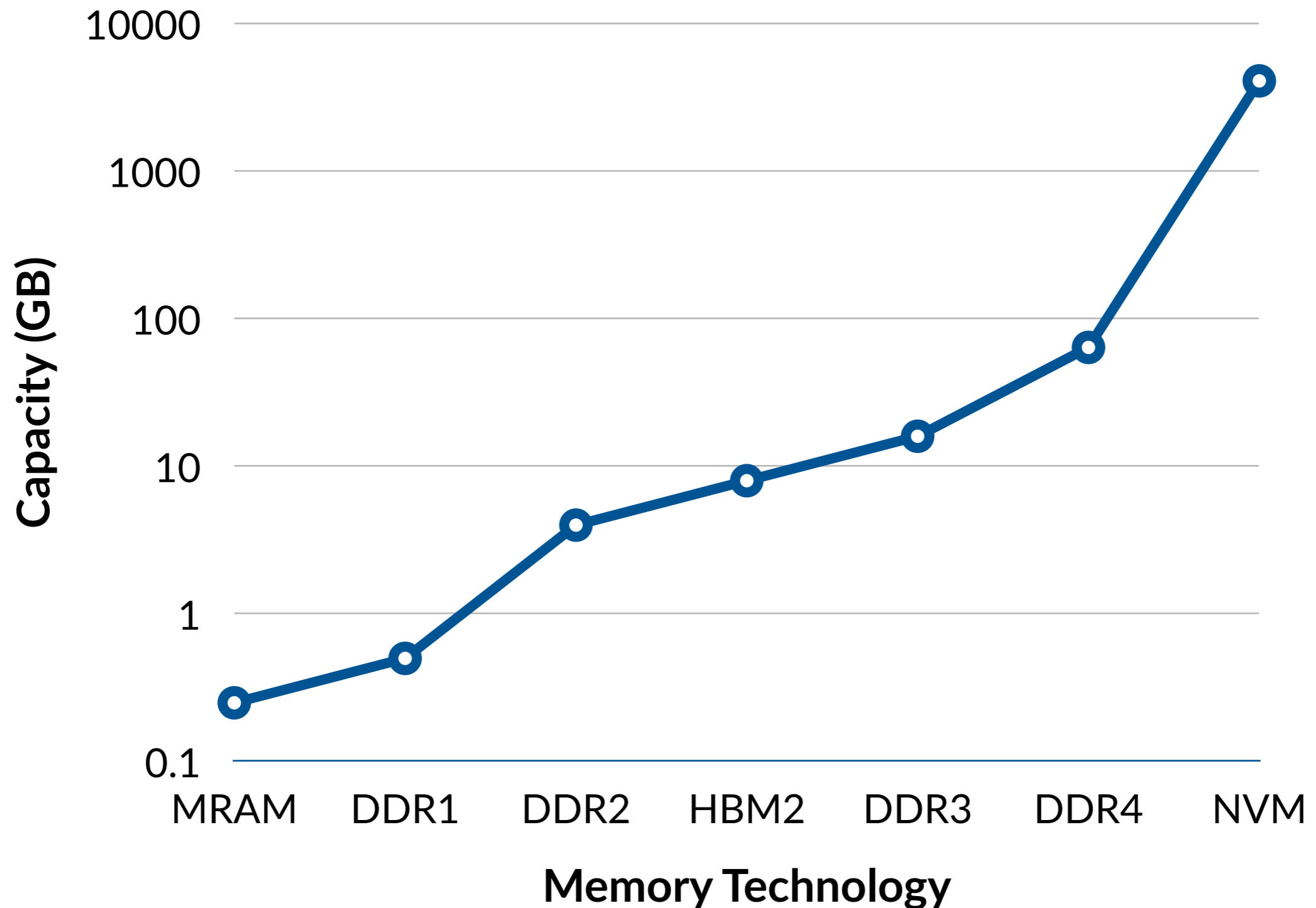
RIP



Memory capacity scaling



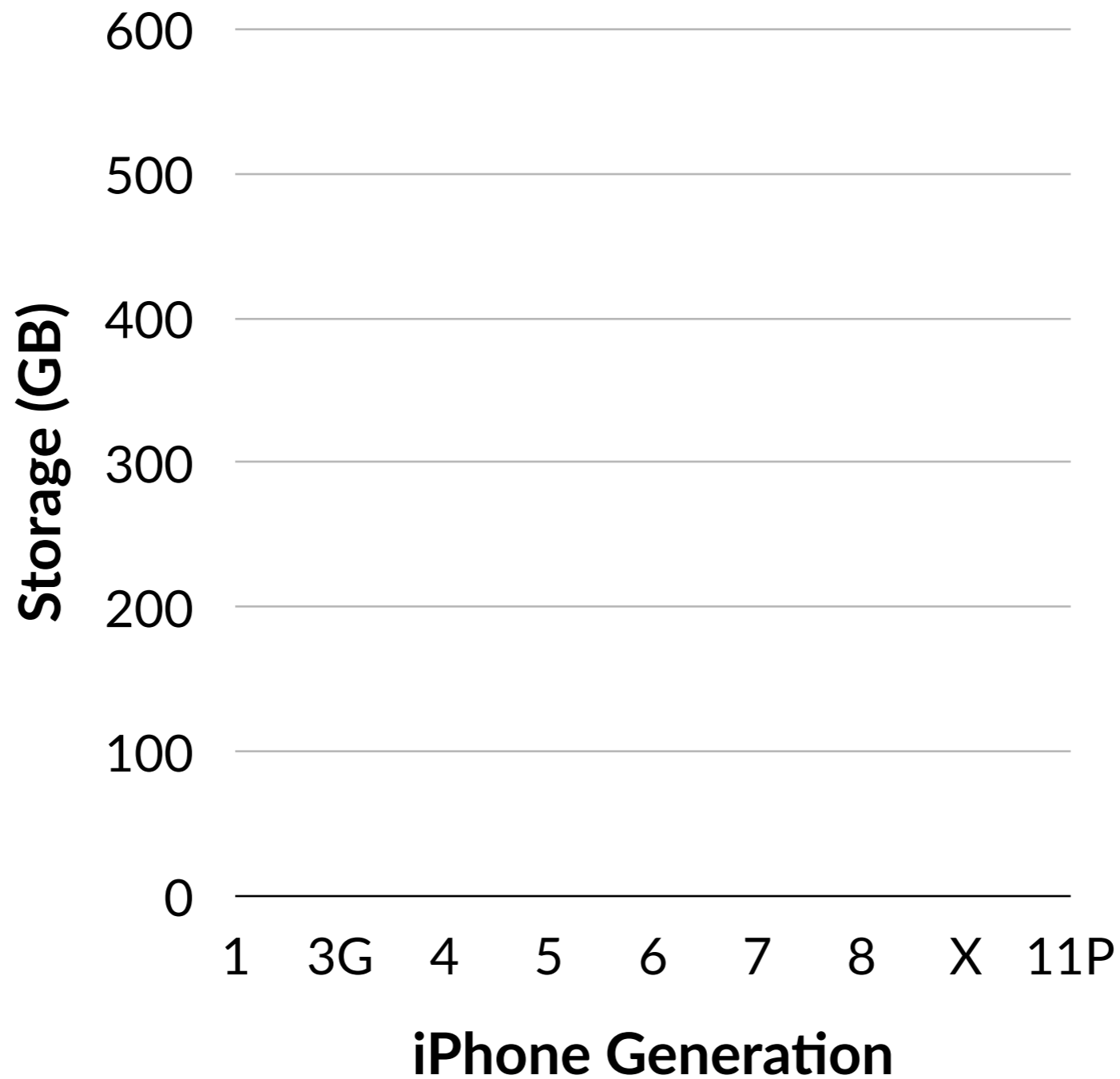
Memory capacity scaling



iPhone storage scaling



2007
iPhone 1
2² GB

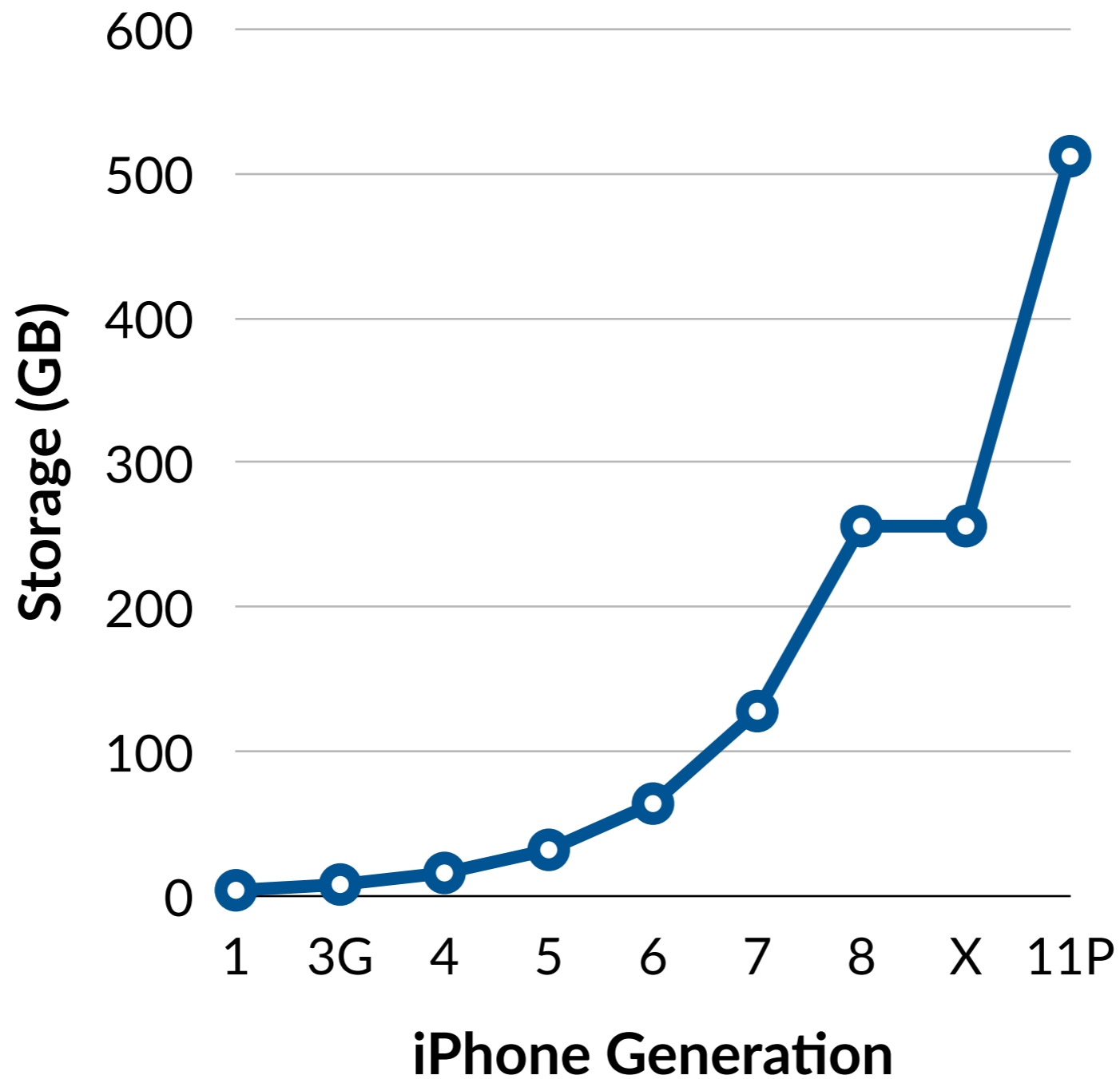


2019
iPhone 11 Pro
2⁹ GB

iPhone storage scaling

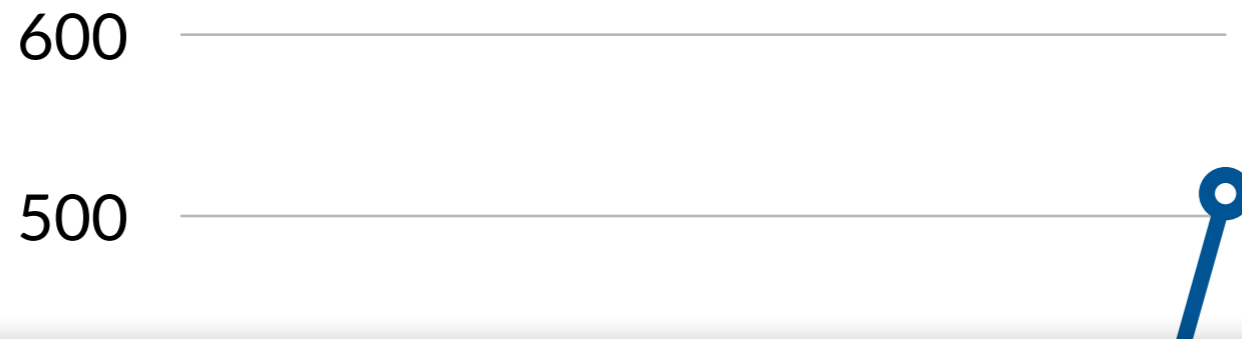
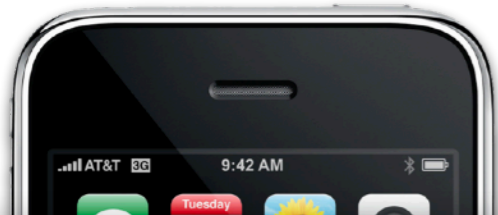


2007
iPhone 1
2² GB



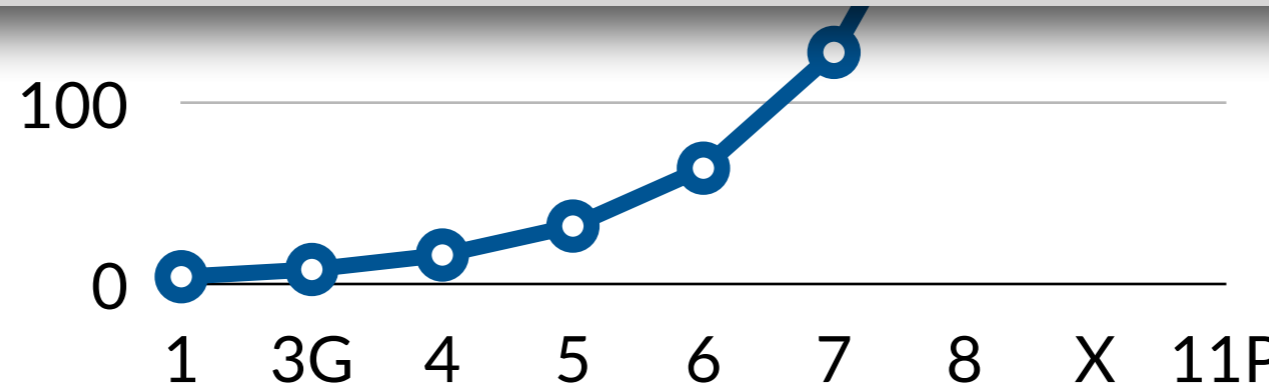
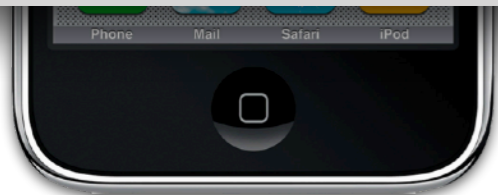
2019
iPhone 11 Pro
2⁹ GB

iPhone storage scaling



At this rate, the 2055 iPhone 47 will store $\sim 2^{75}$ bytes:
the total data *ever generated* thus far

**assumes future iPhone generations are sequential natural numbers*

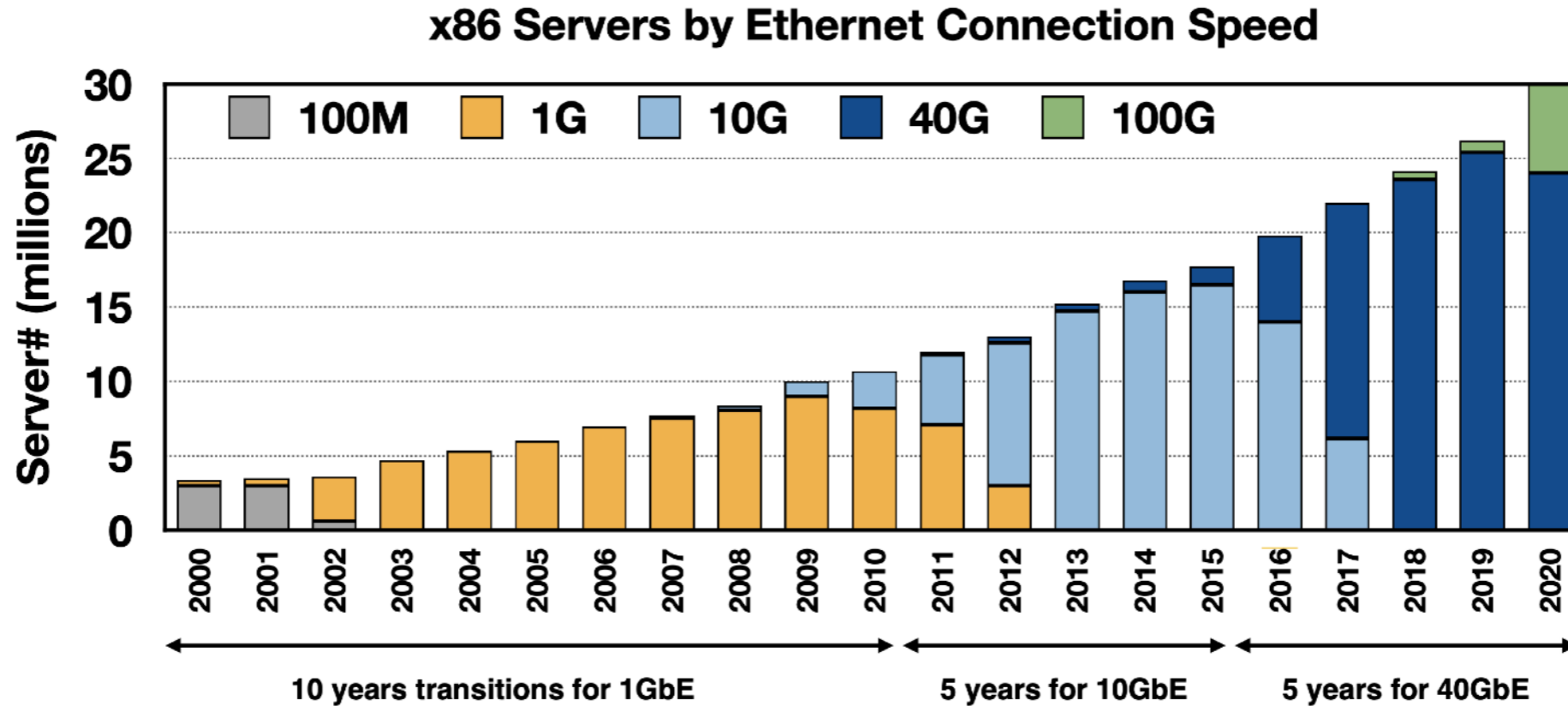


2007
iPhone 1
2² GB

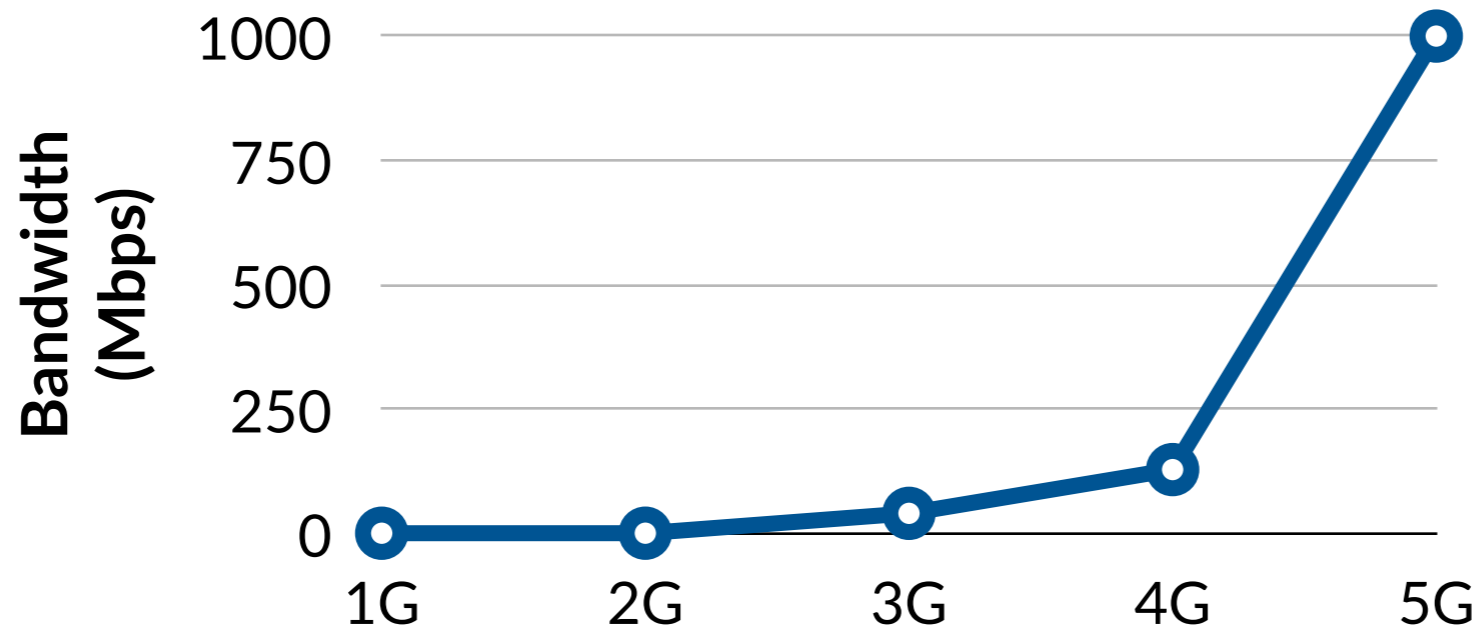
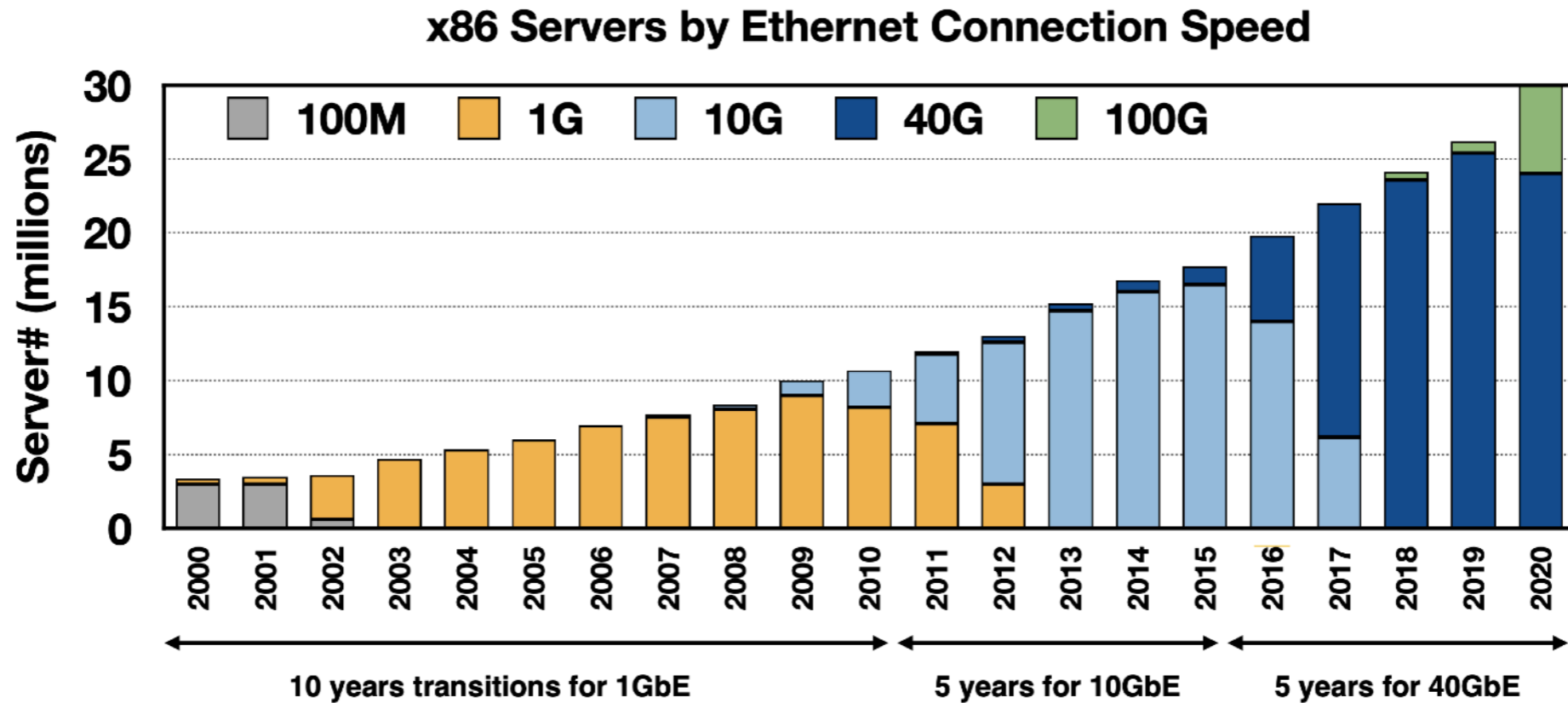
2019
iPhone 11 Pro
2⁹ GB

Datacenter and mobile networks scaling

Datacenter and mobile networks scaling



Datacenter and mobile networks scaling



So let's memoize *almost* everything!



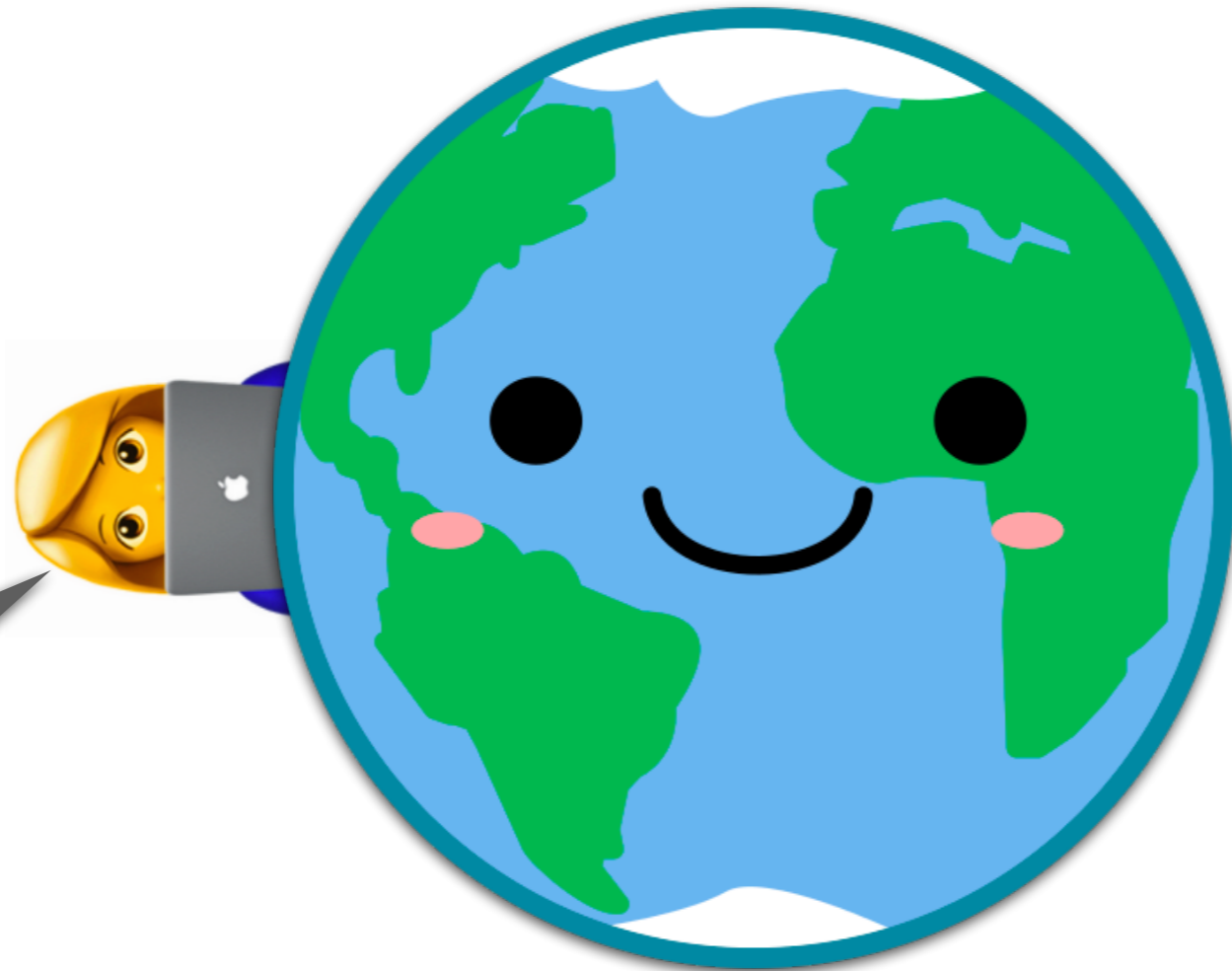
So let's memoize *almost* everything!

Extreme Memoization: store most computation performed and share it globally rather than recomputing!



So let's memoize *almost* everything!

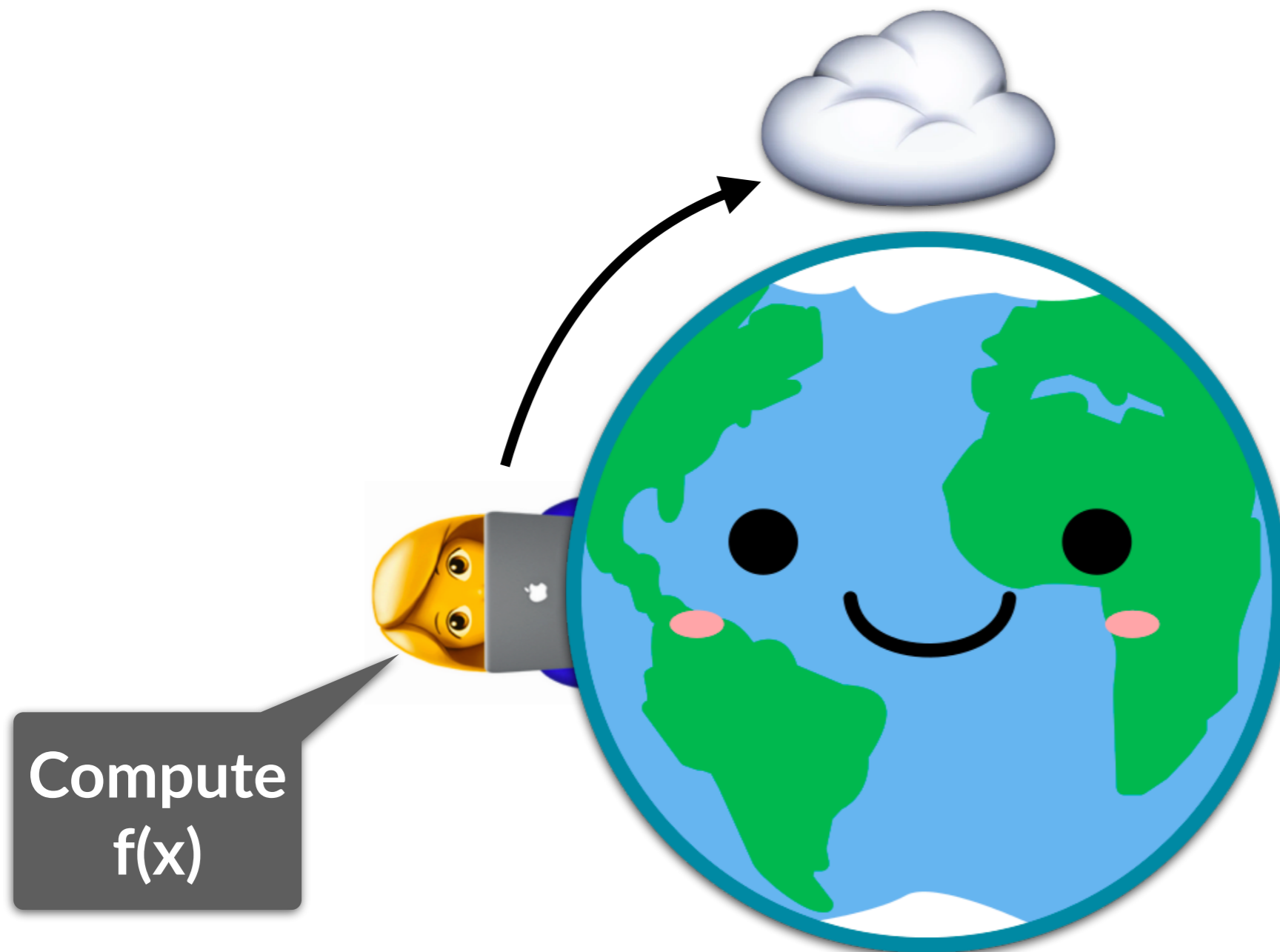
Extreme Memoization: store most computation performed and share it globally rather than recomputing!



Compute
 $f(x)$

So let's memoize *almost* everything!

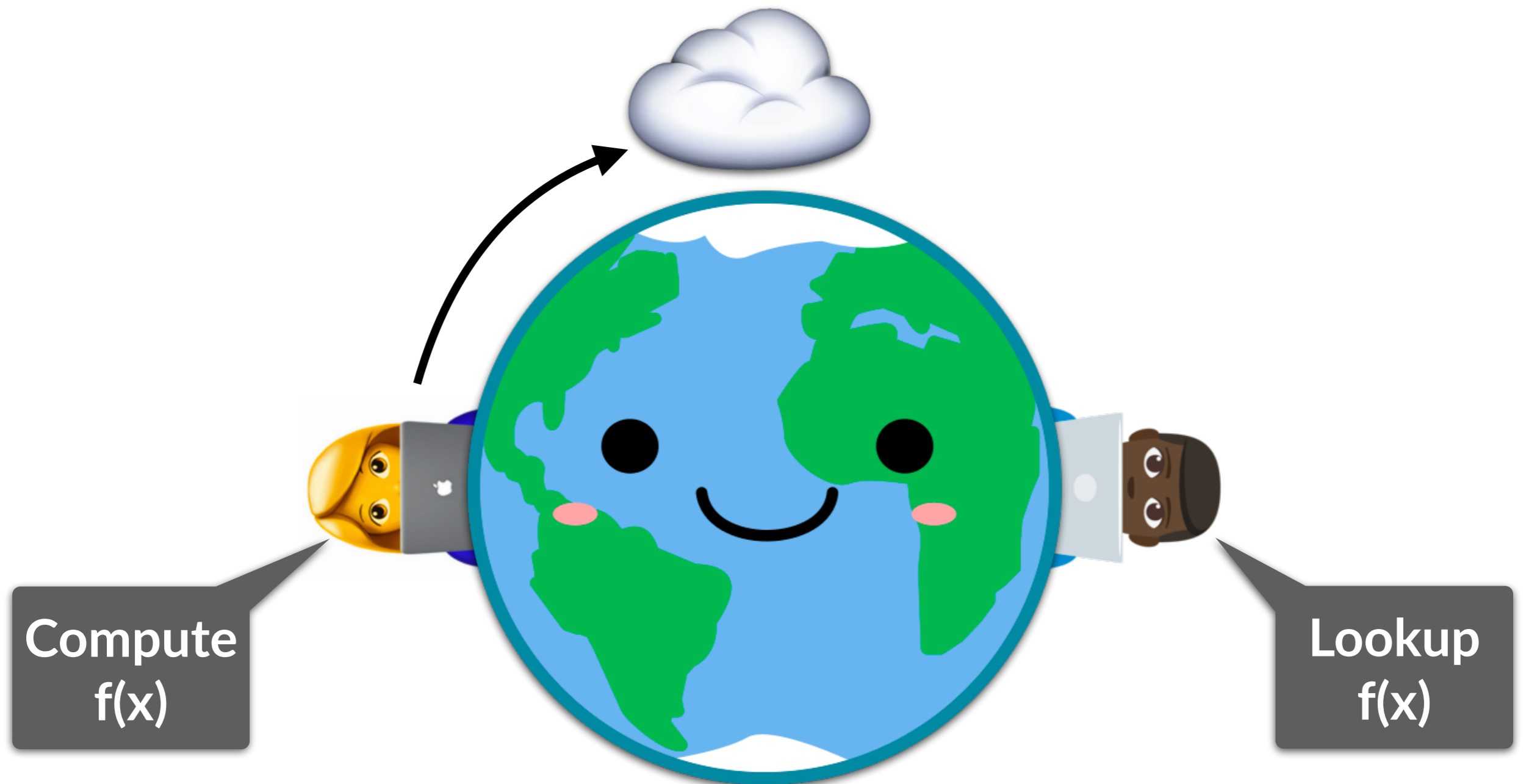
Extreme Memoization: store most computation performed and share it globally rather than recomputing!



Compute
 $f(x)$

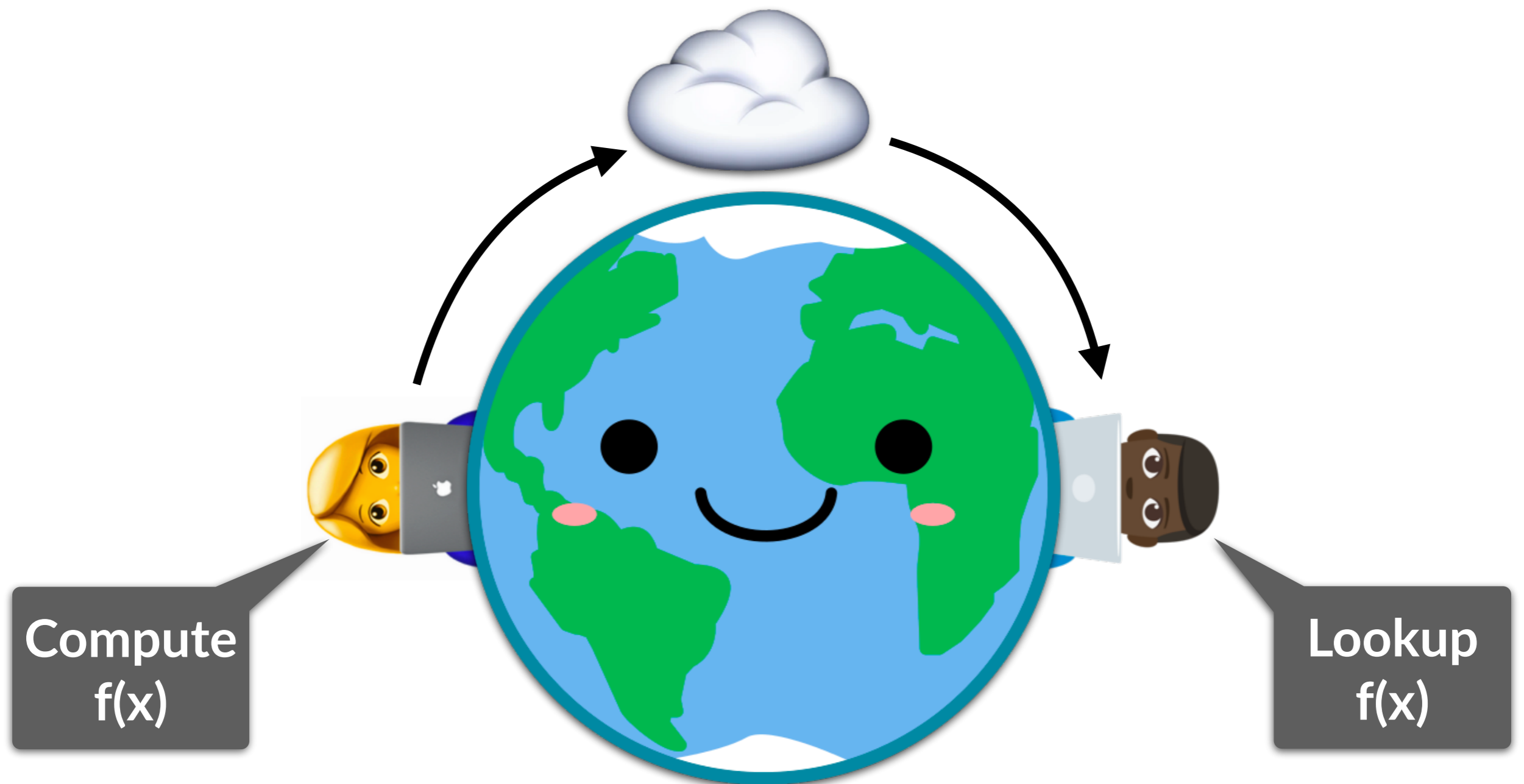
So let's memoize *almost* everything!

Extreme Memoization: store most computation performed and share it globally rather than recomputing!



So let's memoize *almost* everything!

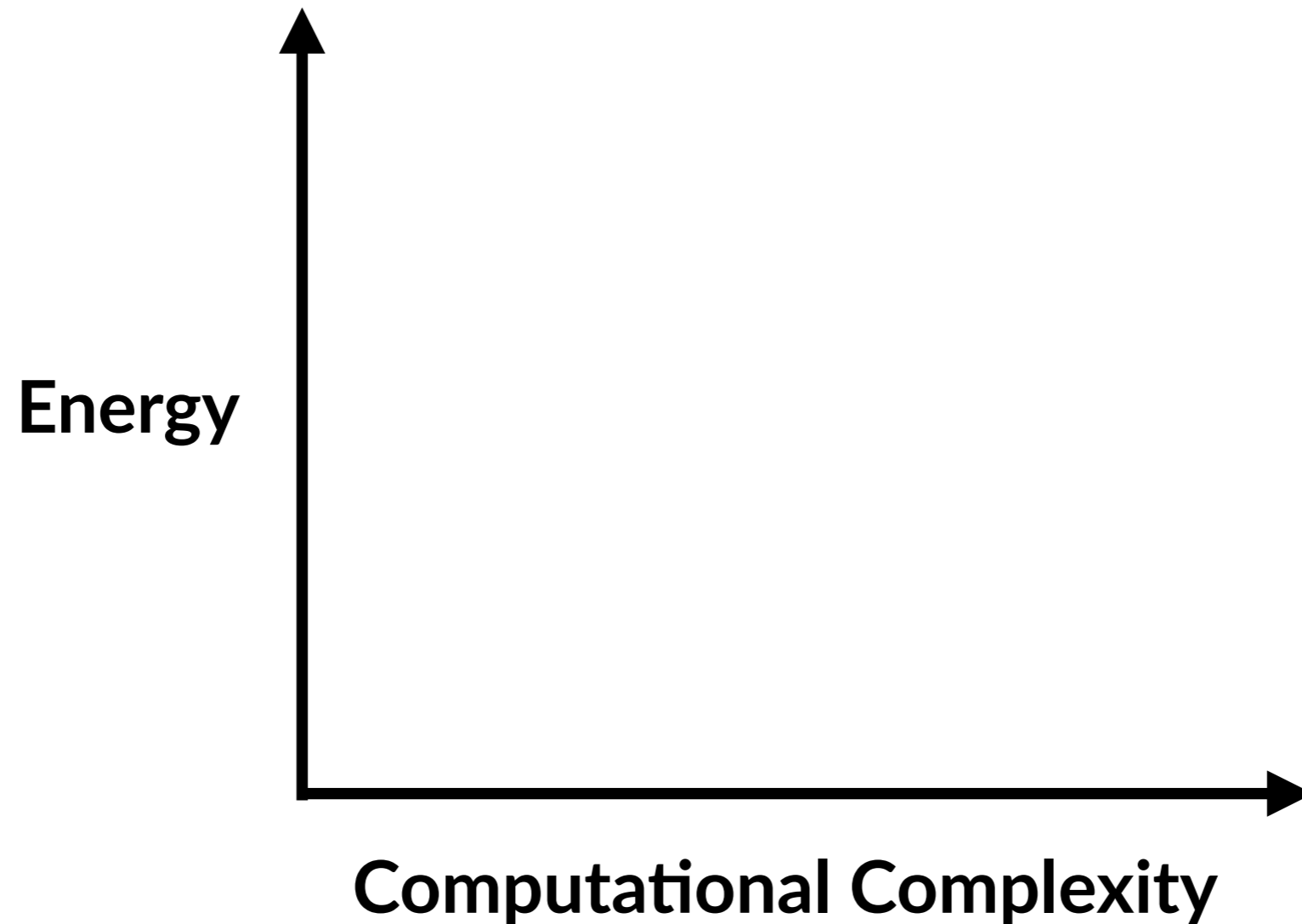
Extreme Memoization: store most computation performed and share it globally rather than recomputing!



When will this make sense?

— CPU compute energy

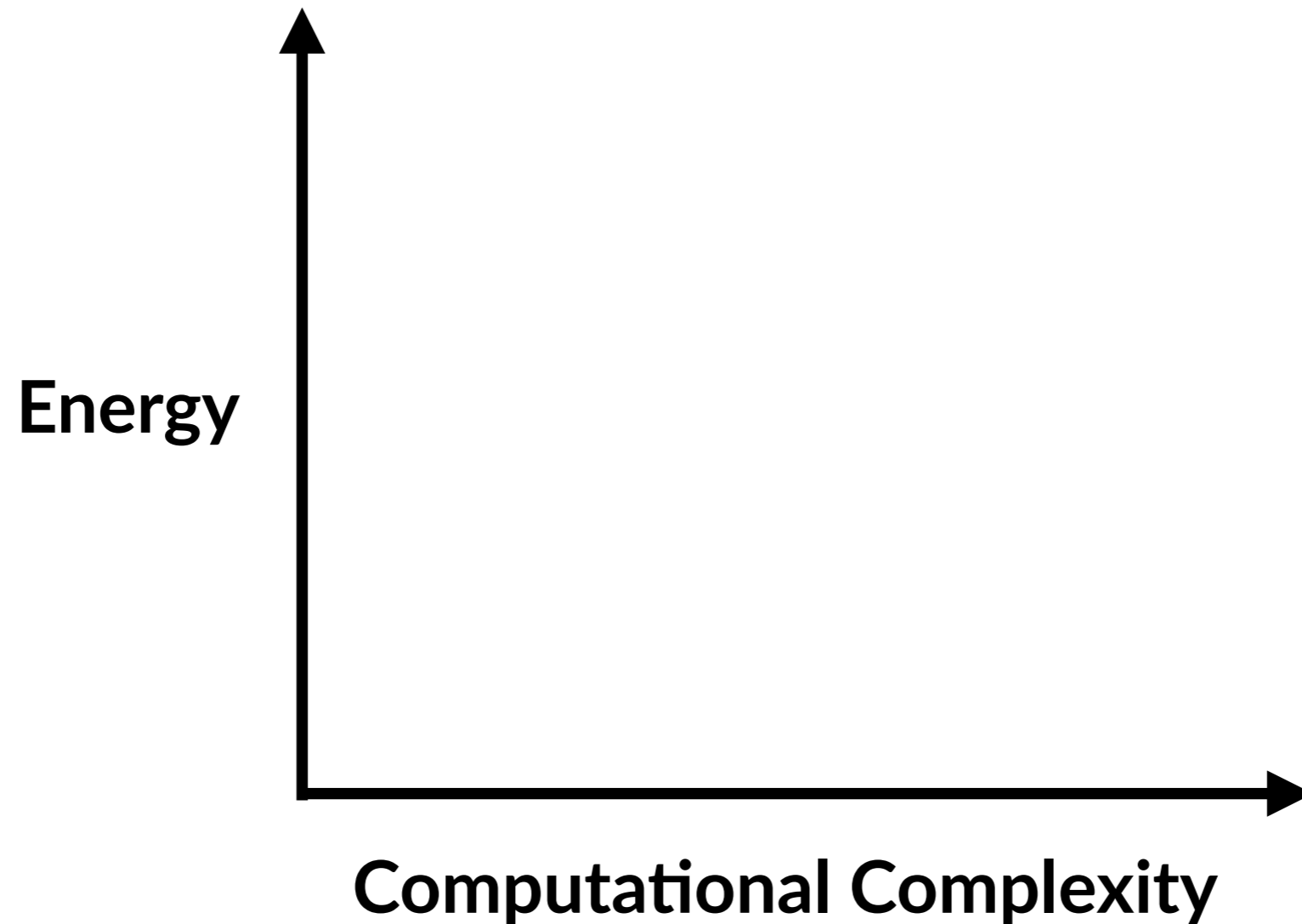
— Memoization energy



When will this make sense?

— CPU compute energy

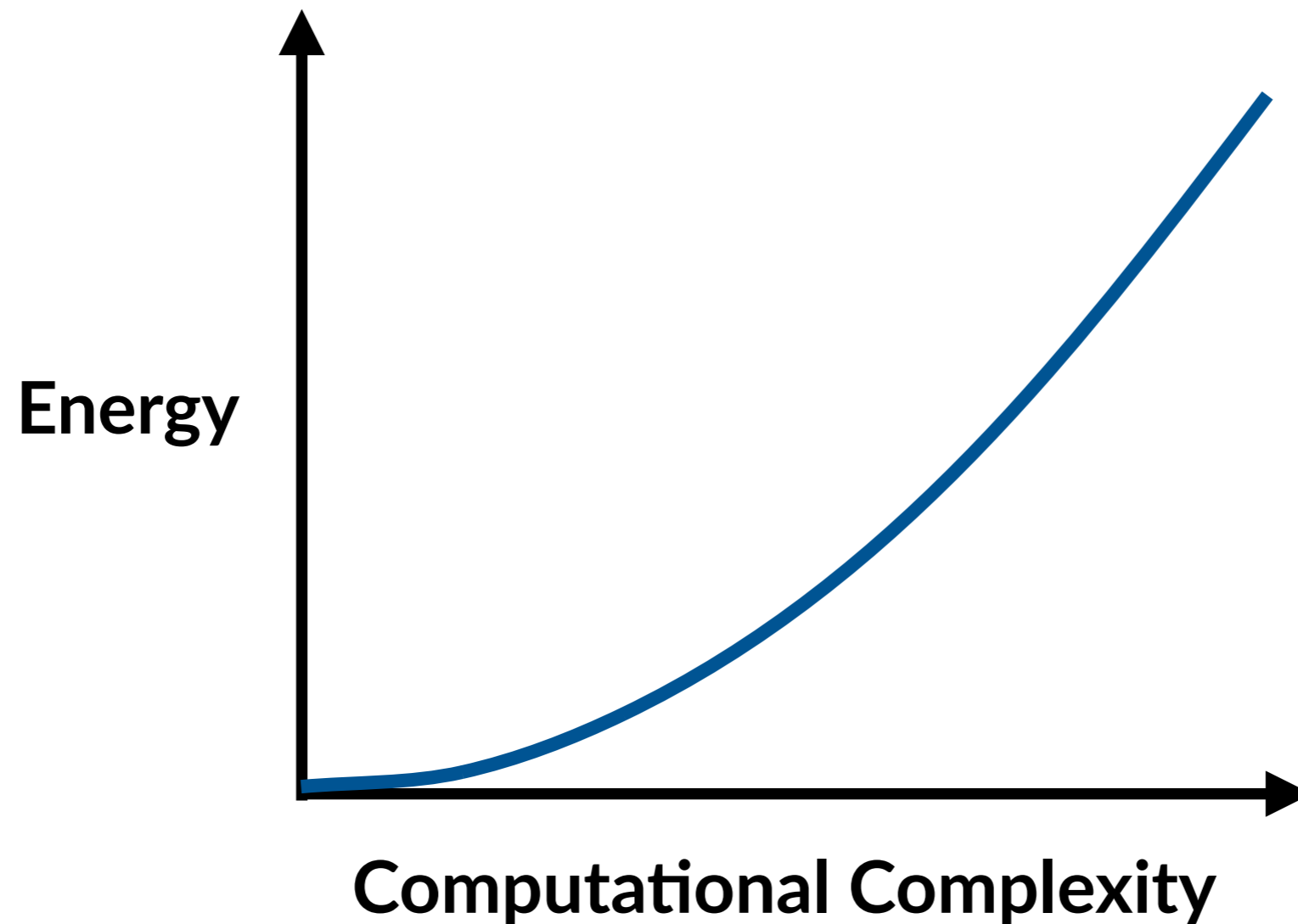
— Memoization energy



When will this make sense?

— CPU compute energy

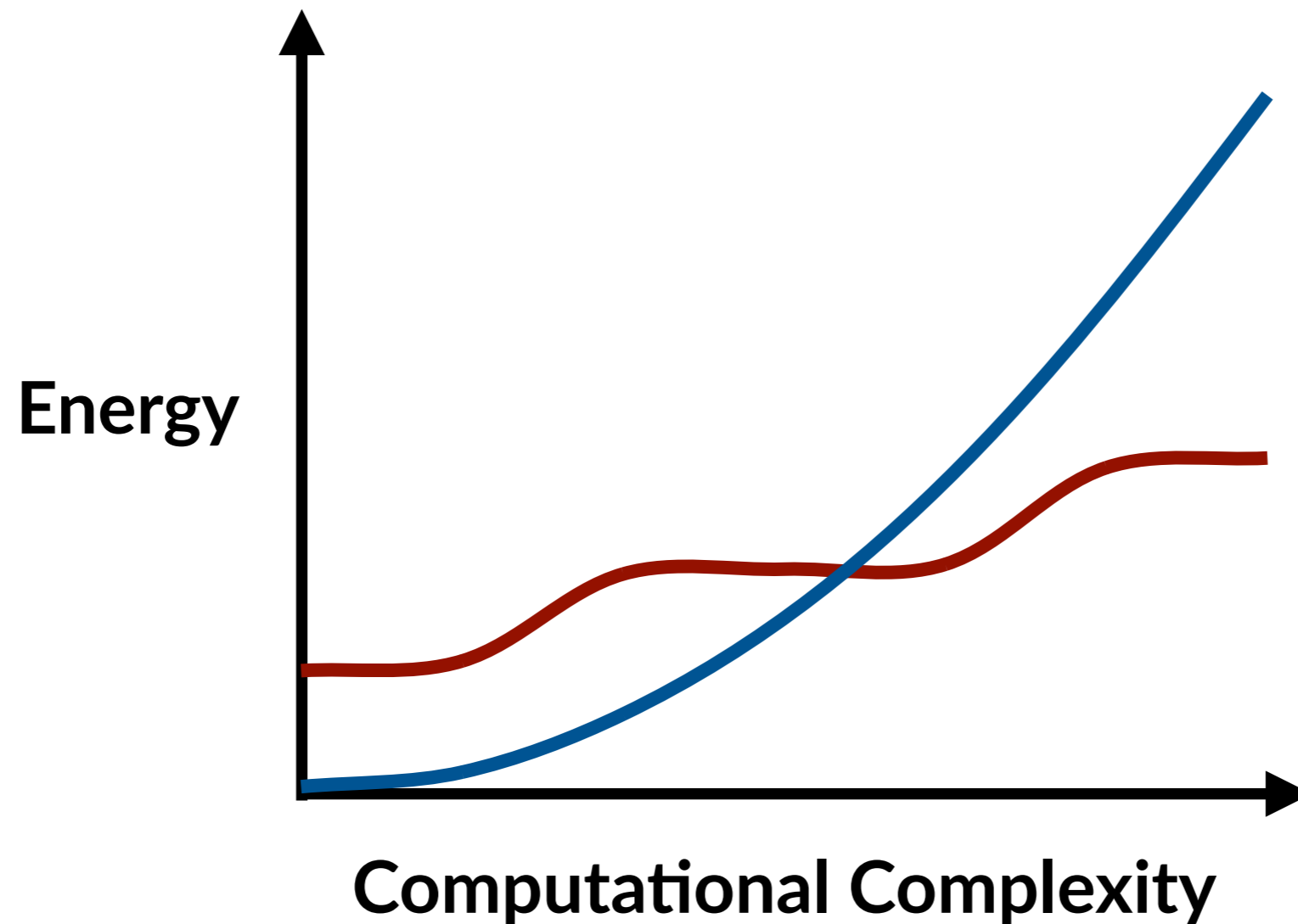
— Memoization energy



When will this make sense?

— CPU compute energy

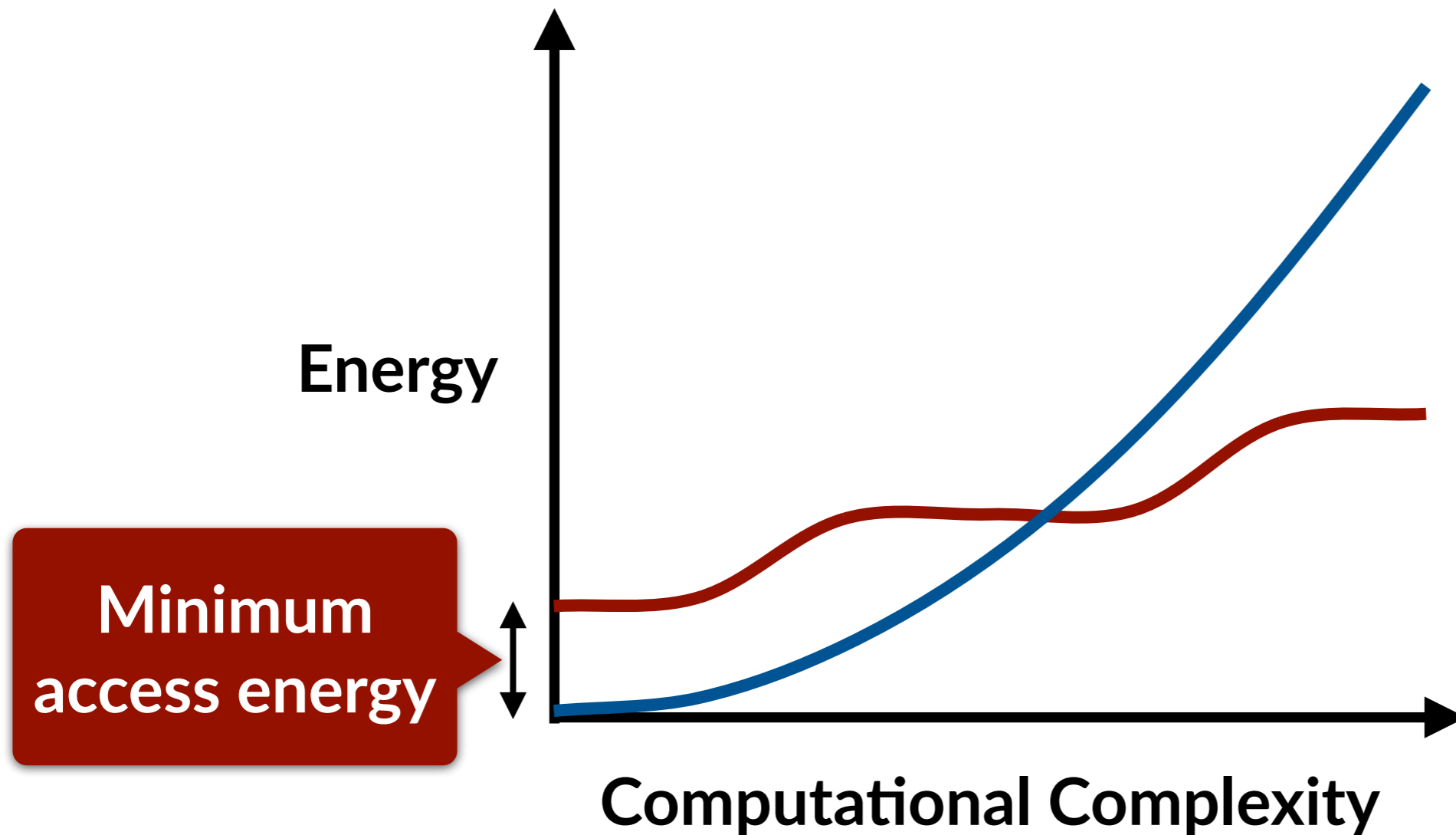
— Memoization energy



When will this make sense?

— CPU compute energy

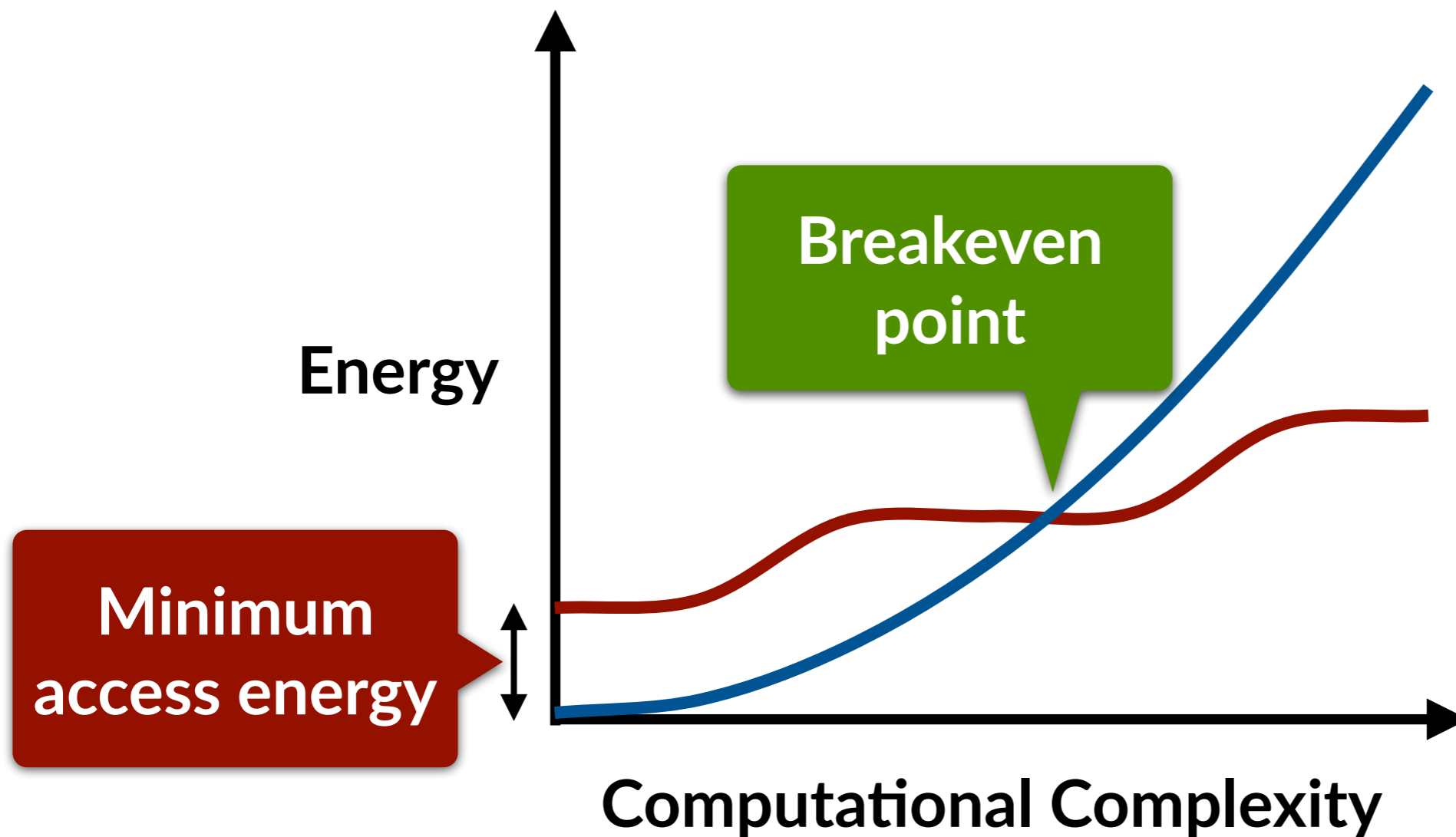
— Memoization energy



When will this make sense?

— CPU compute energy

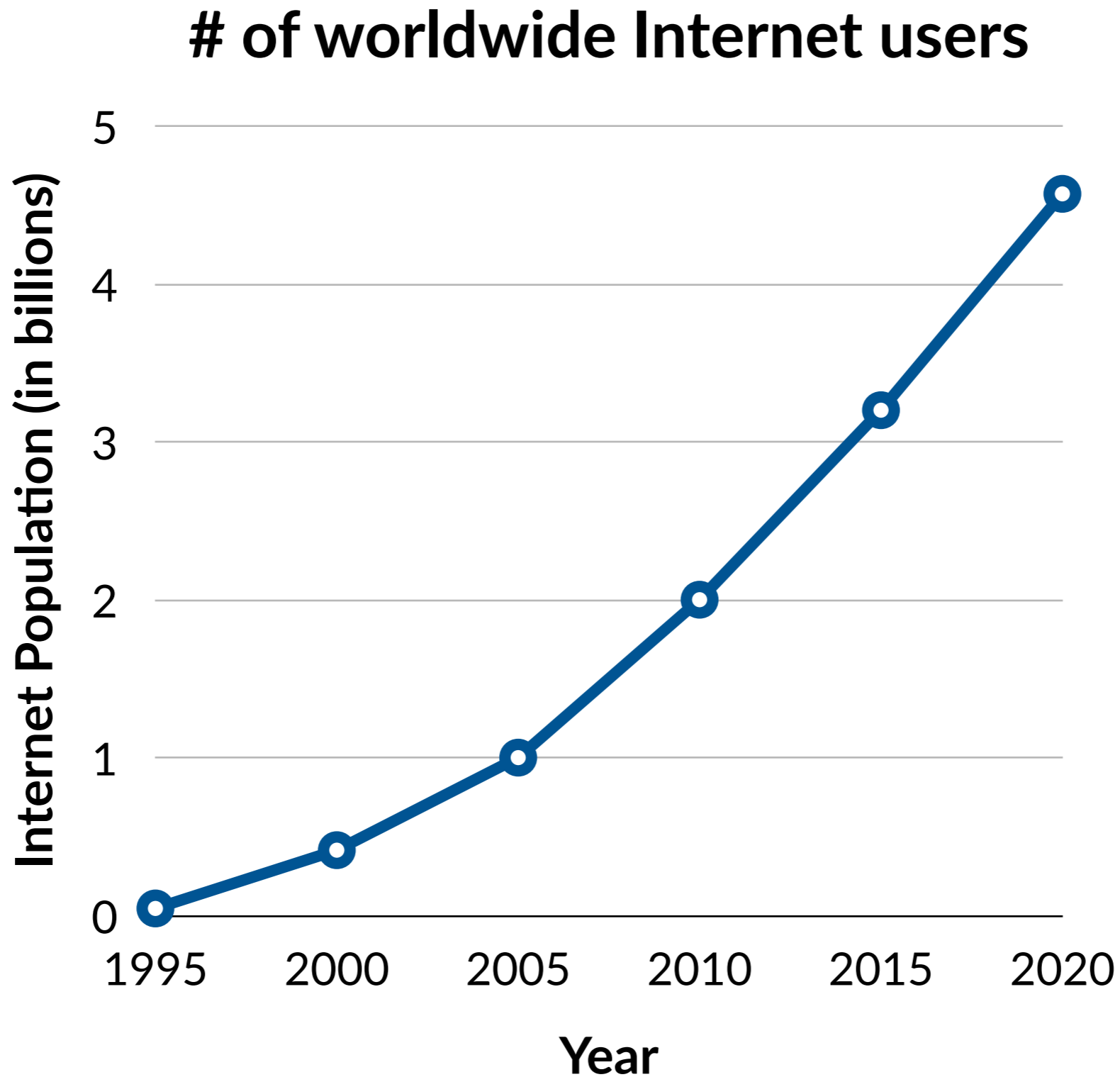
— Memoization energy



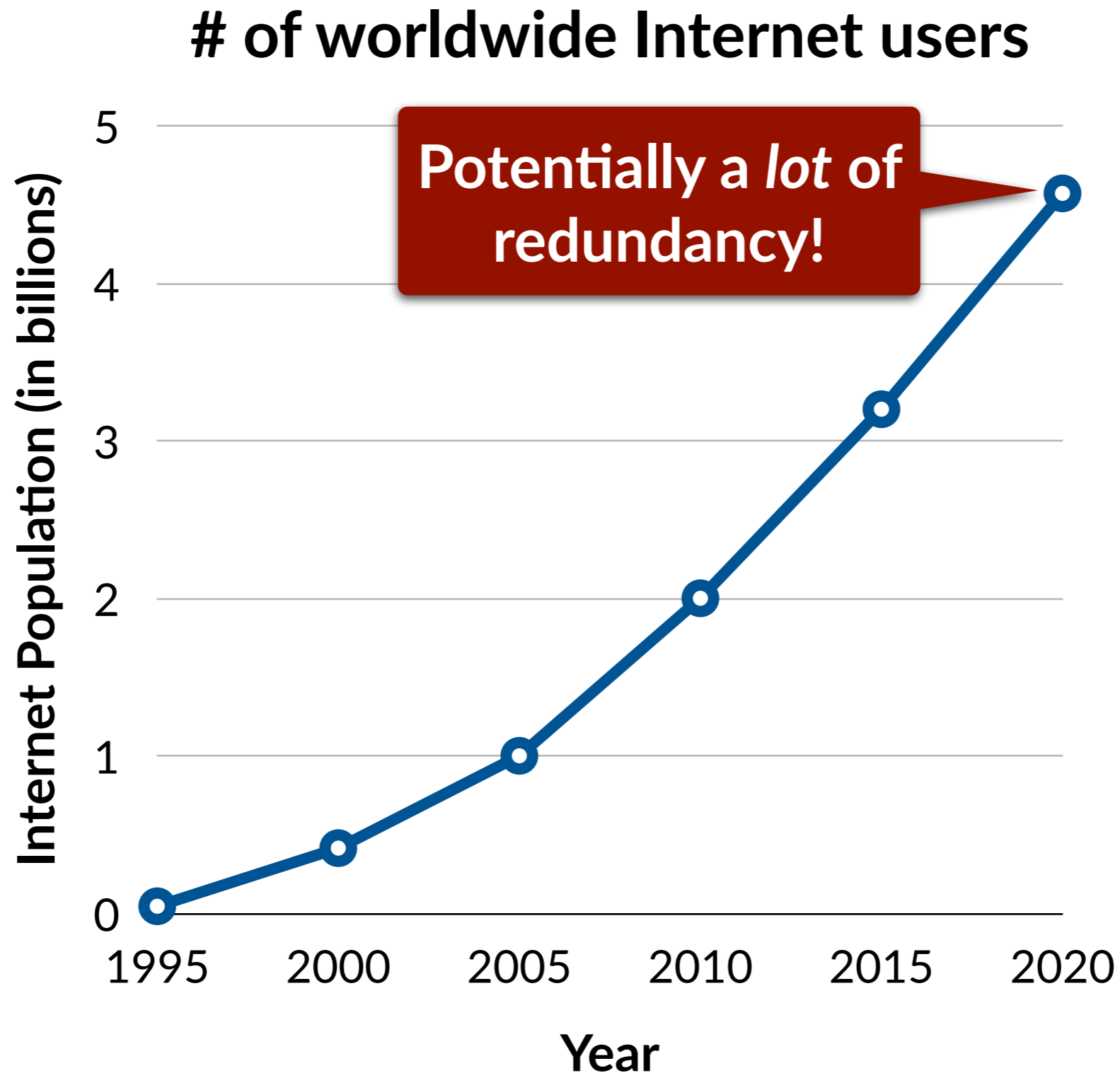
When storage + access energy is lower than computation energy

Computation redundancy wastes energy

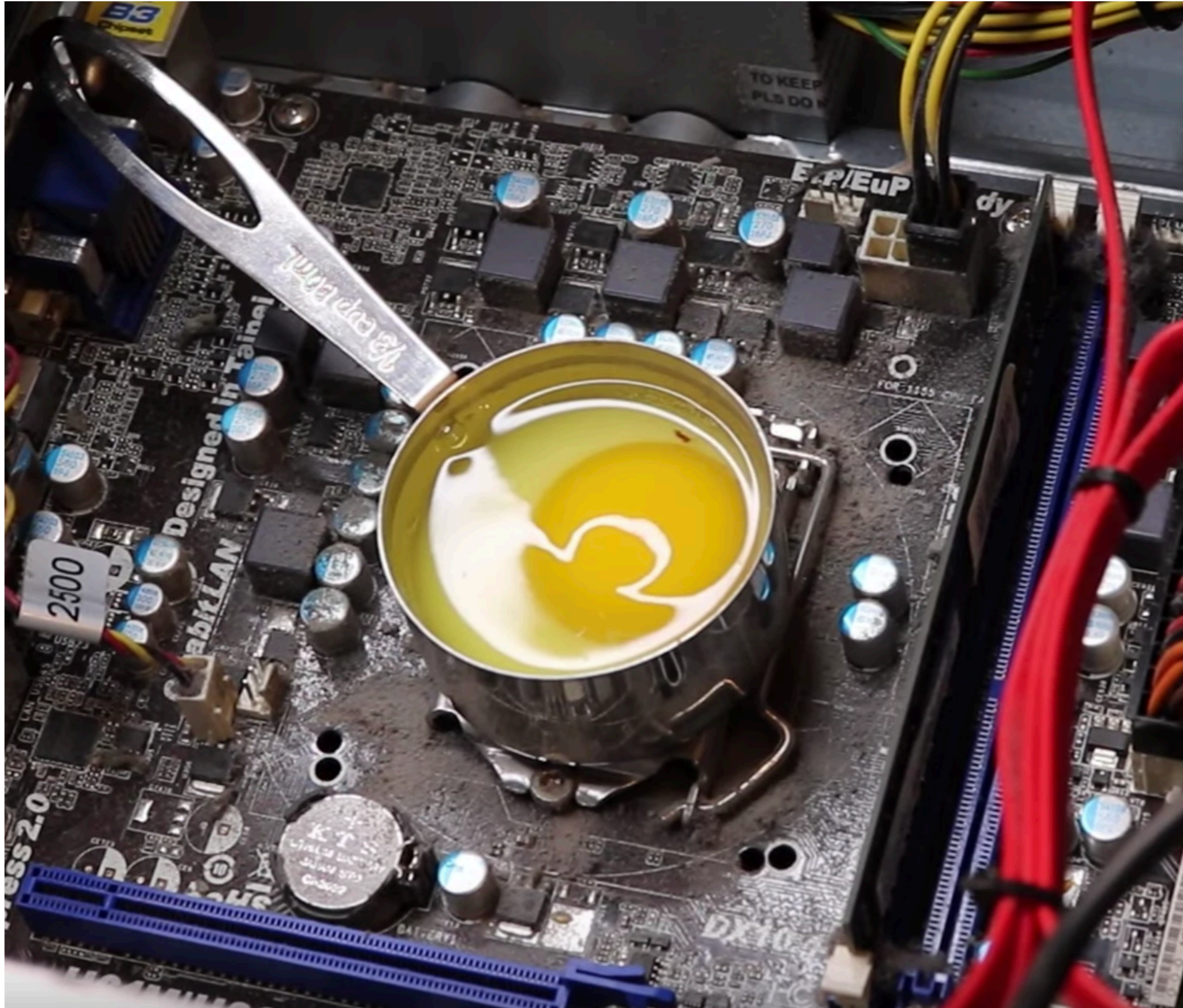
Computation redundancy wastes energy



Computation redundancy wastes energy



Computation redundancy wastes energy



We compute by memorization!! 🤔



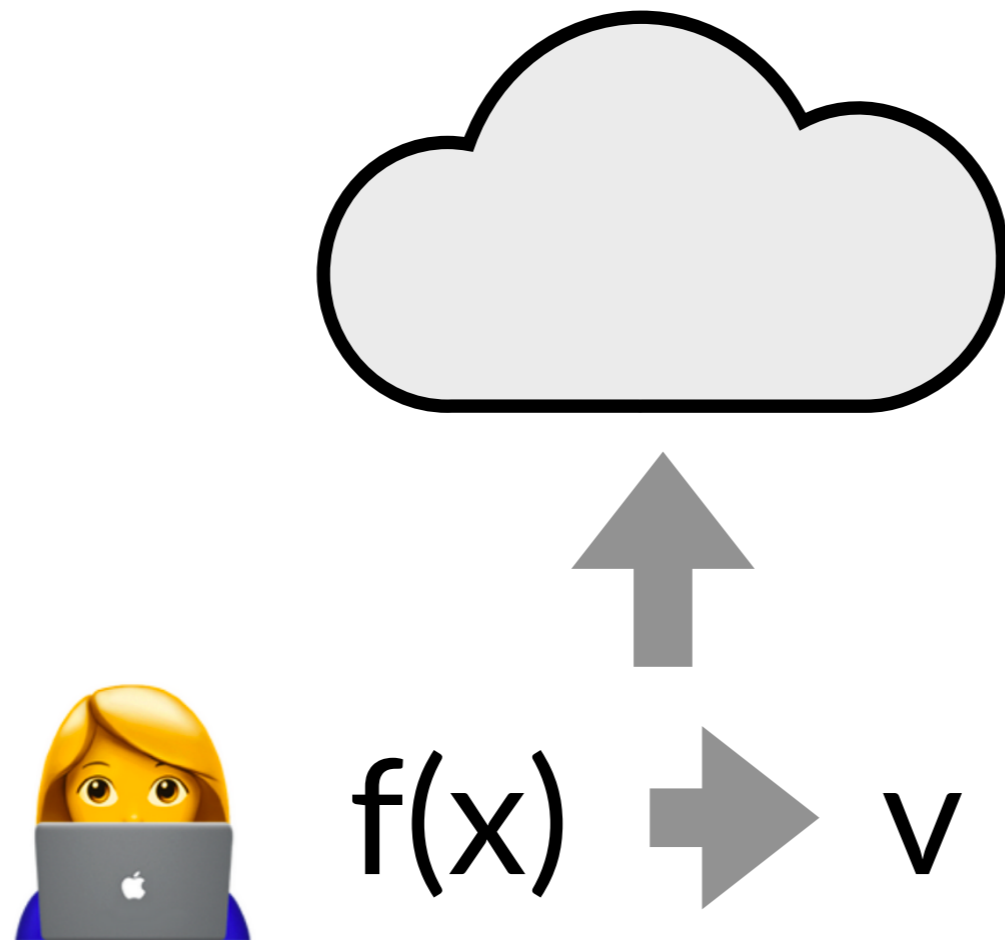
Source: *Harvard University Youtube Channel*

We compute by memorization!! 🧠

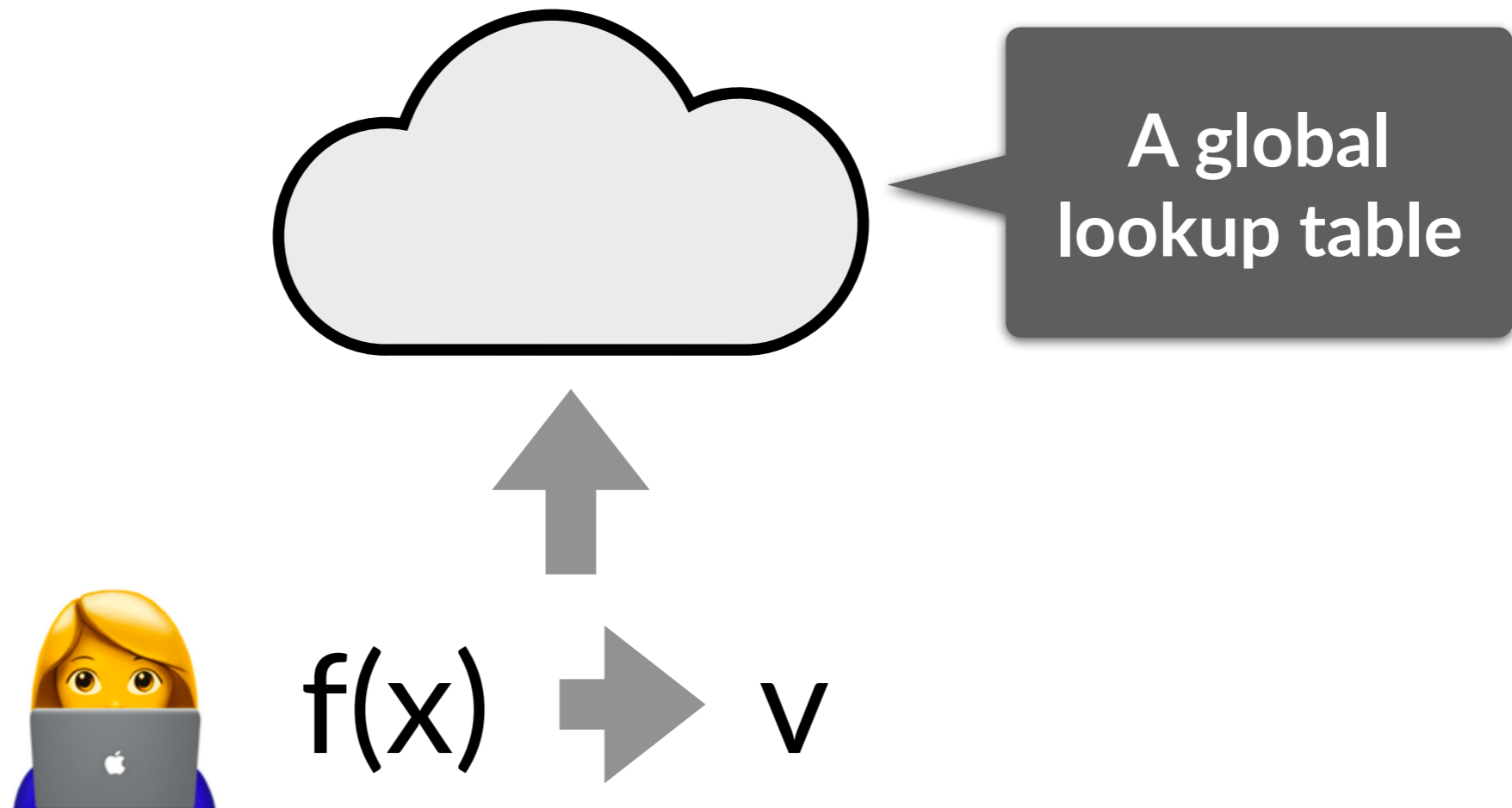


Source: *Harvard University Youtube Channel*

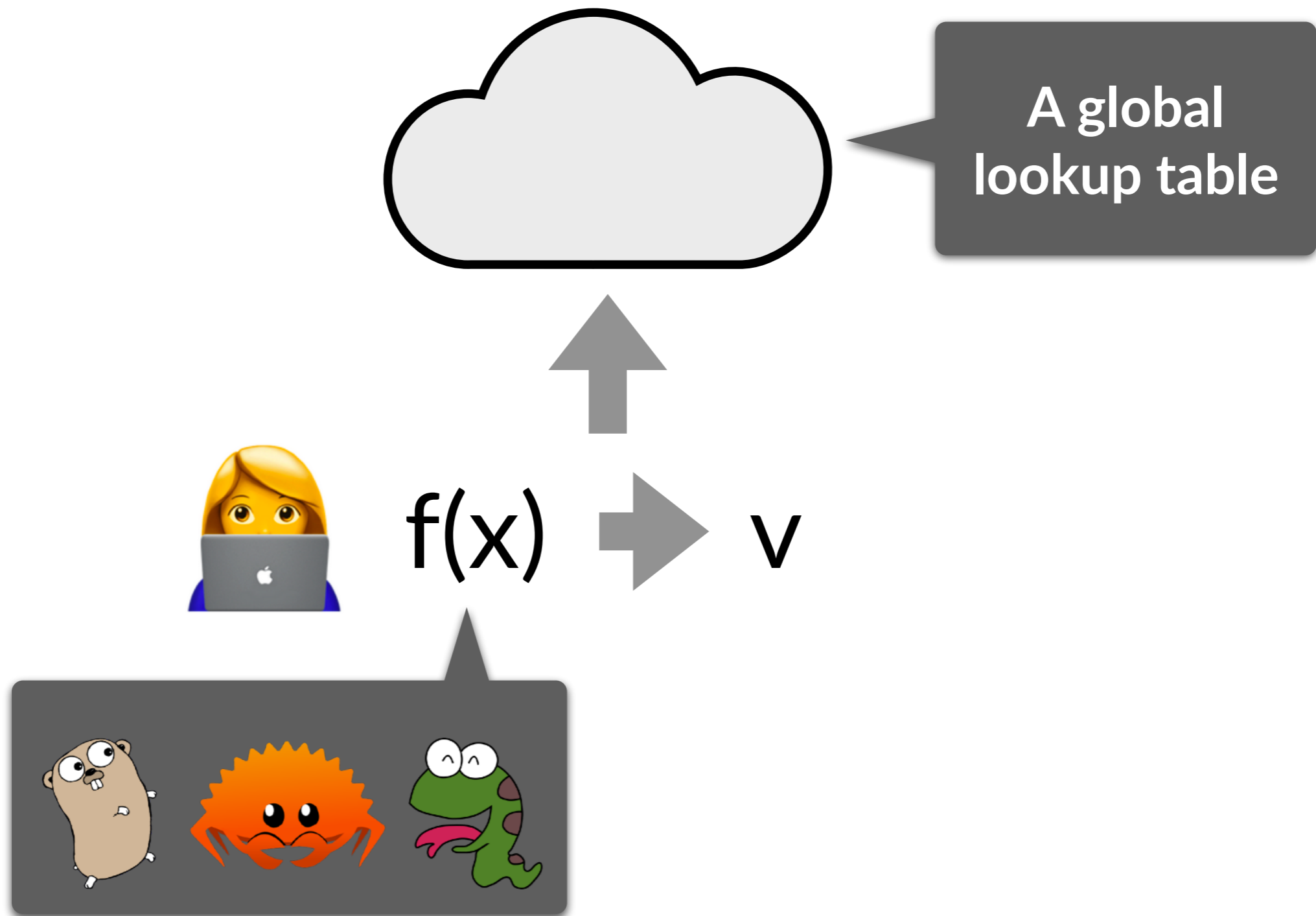
Memoize function, inputs, and outputs



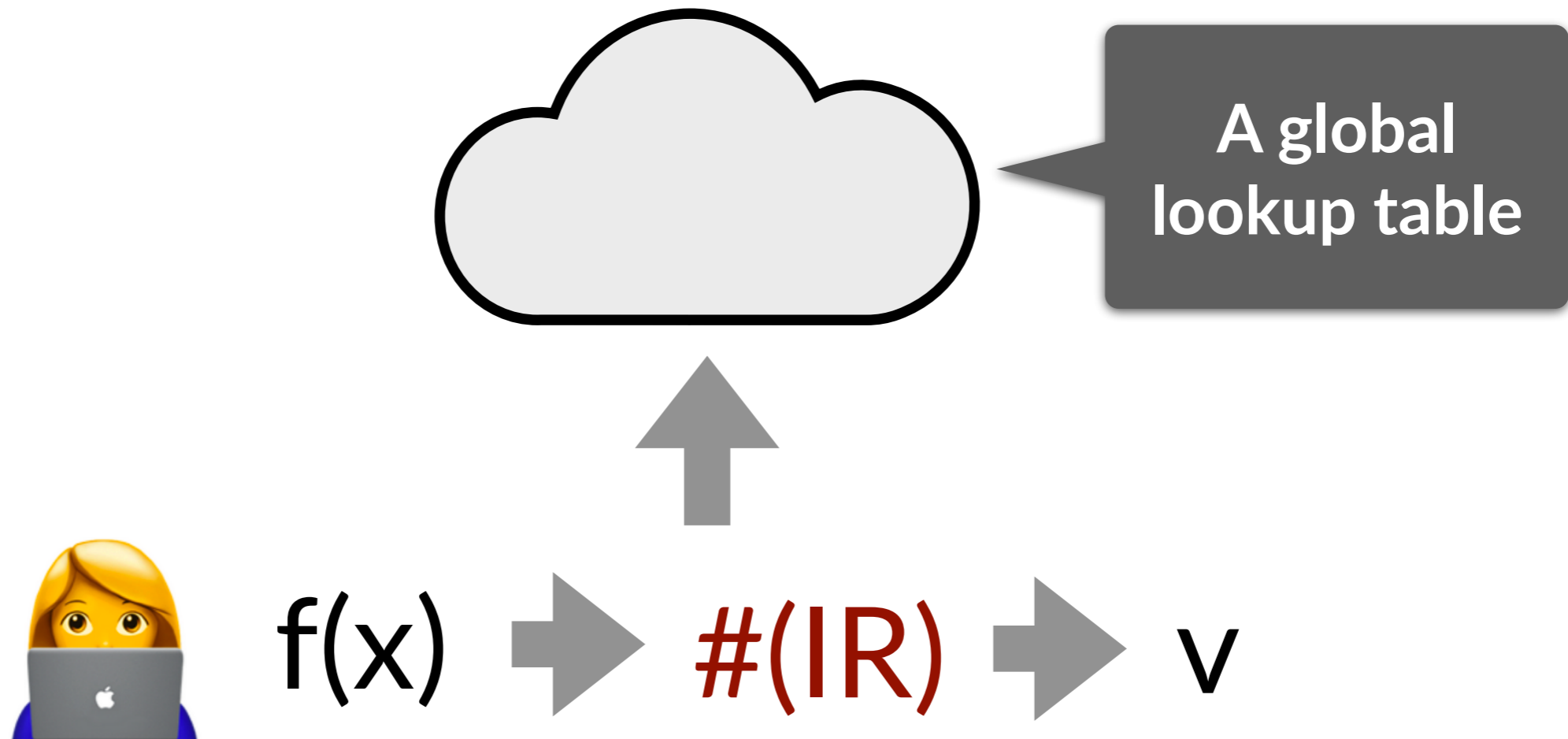
Memoize function, inputs, and outputs



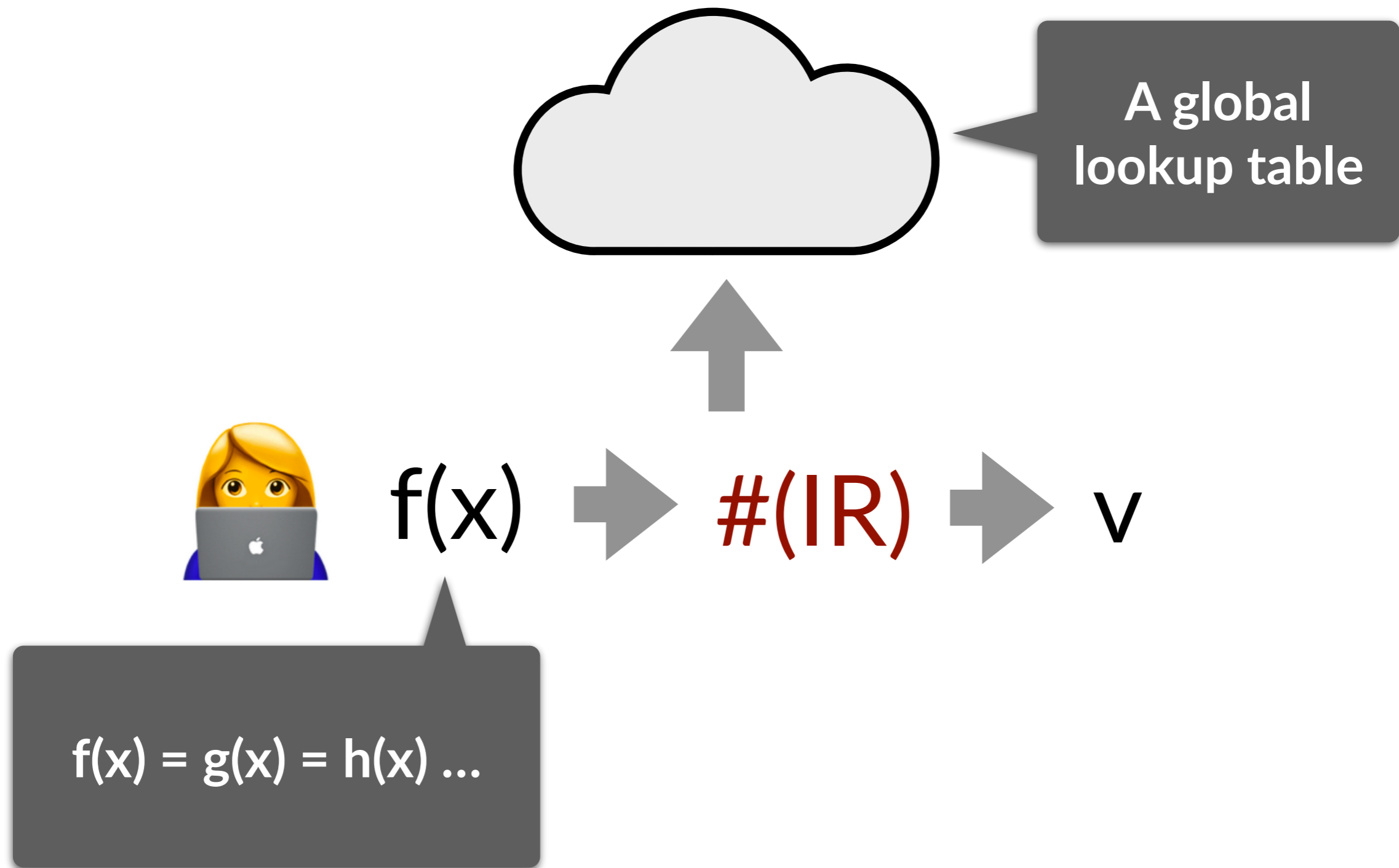
f() implemented in any PL



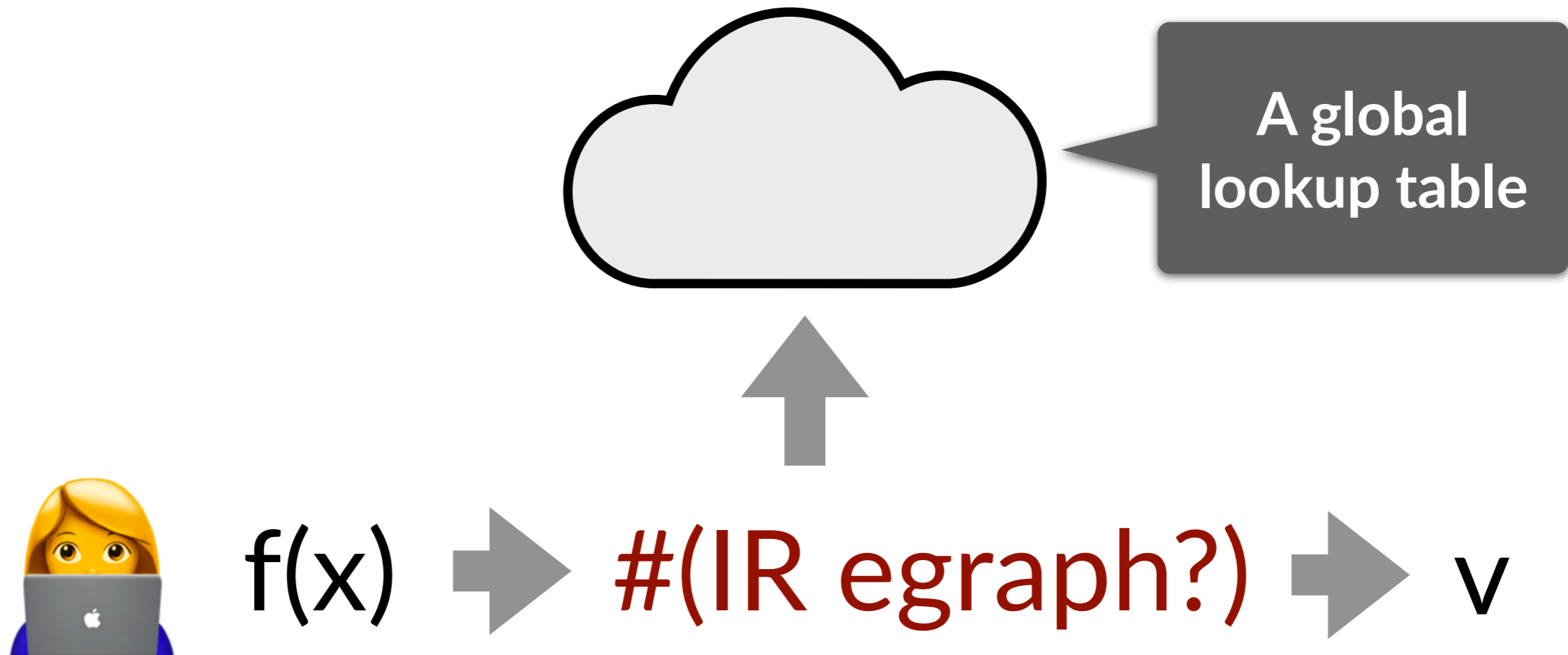
Hash a language-agnostic IR!



Many ways to write the same function



Maybe use equivalence graphs?



Maybe use equivalence graphs?

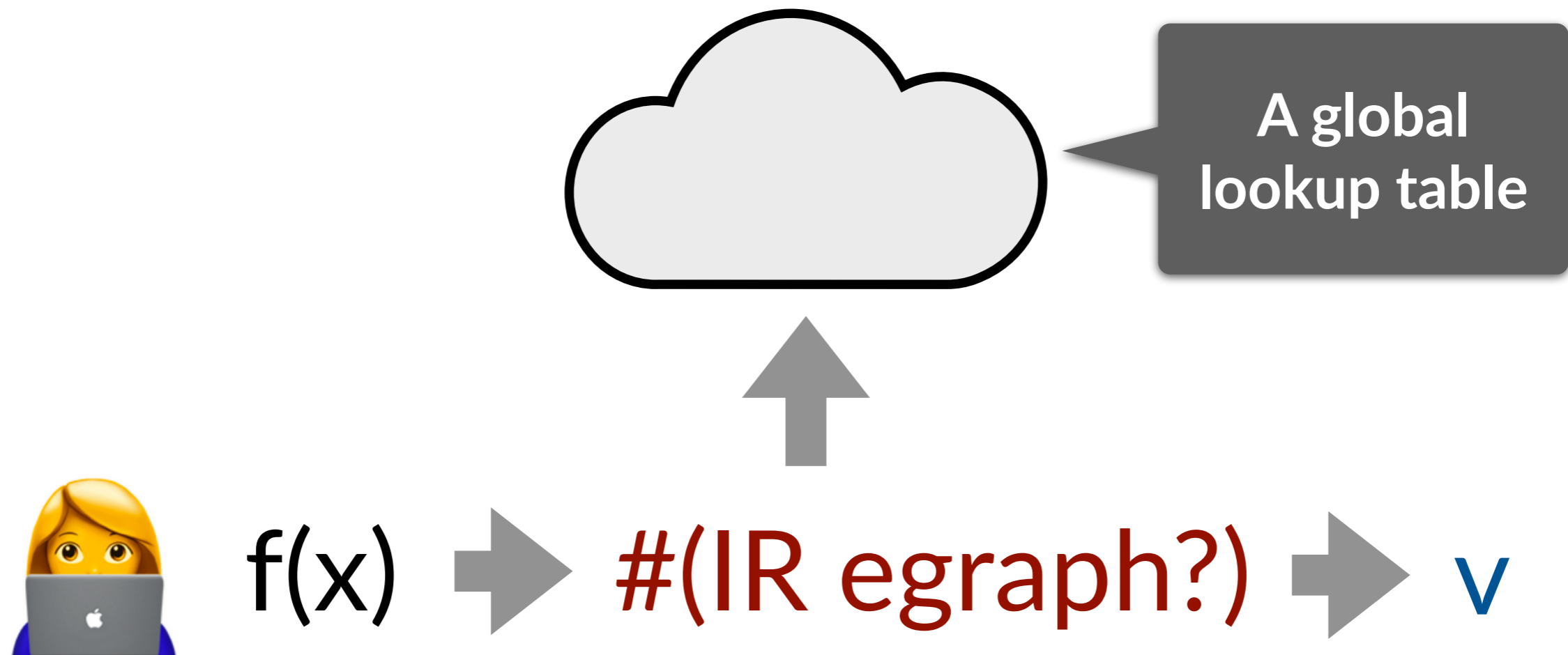


Goal
table

v



Store **key**-**value** pair in the global LUT



Too big a (lookup) table?



Too big a (lookup) table?

Worldwide compute capacity
is $\sim 10^{20}$ FLOPS!

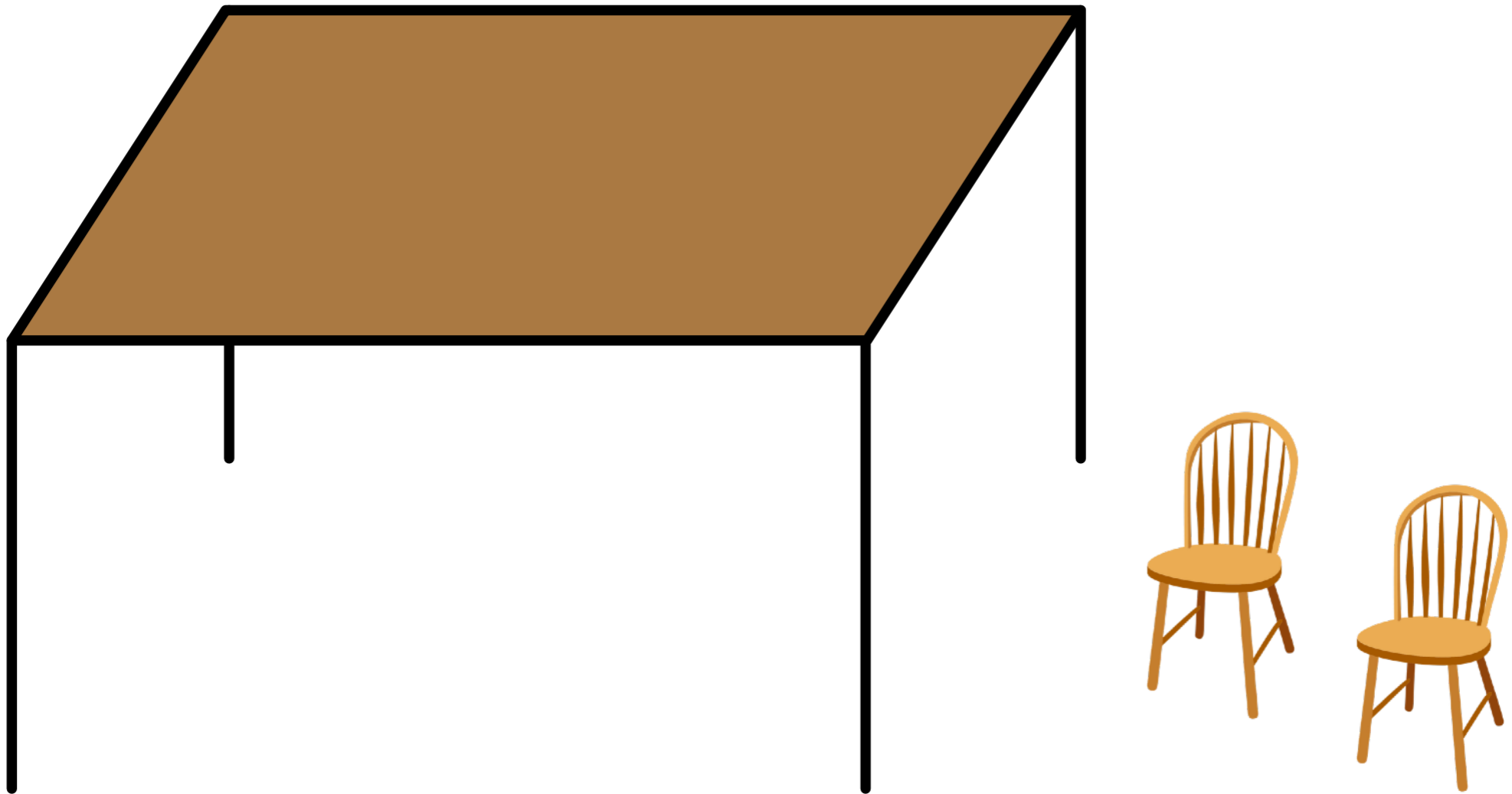


We might need a helluva lot of Hellabytes

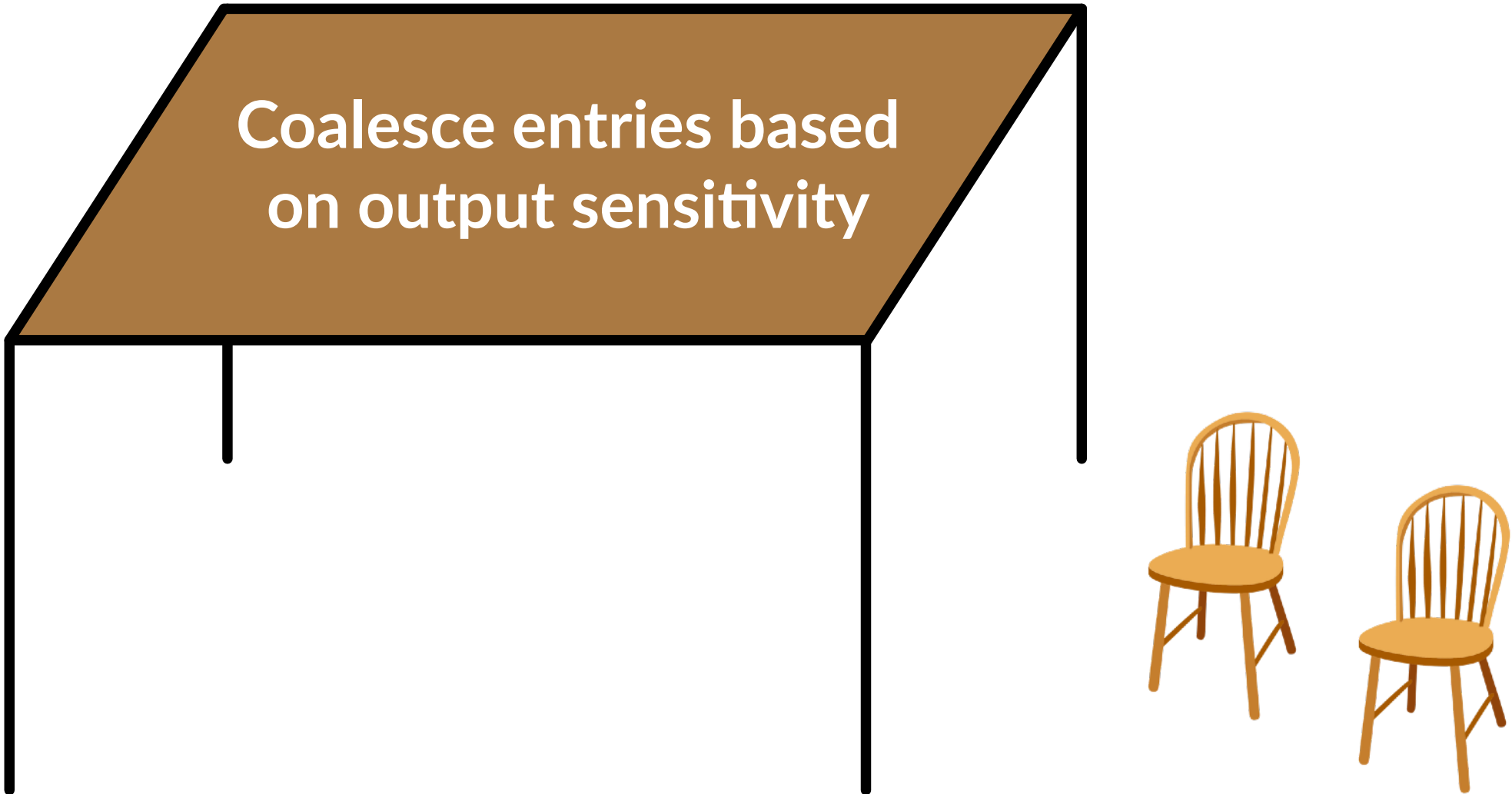
1 Hellabyte = 10^{27} bytes



Approximation can help...



Approximation can help...



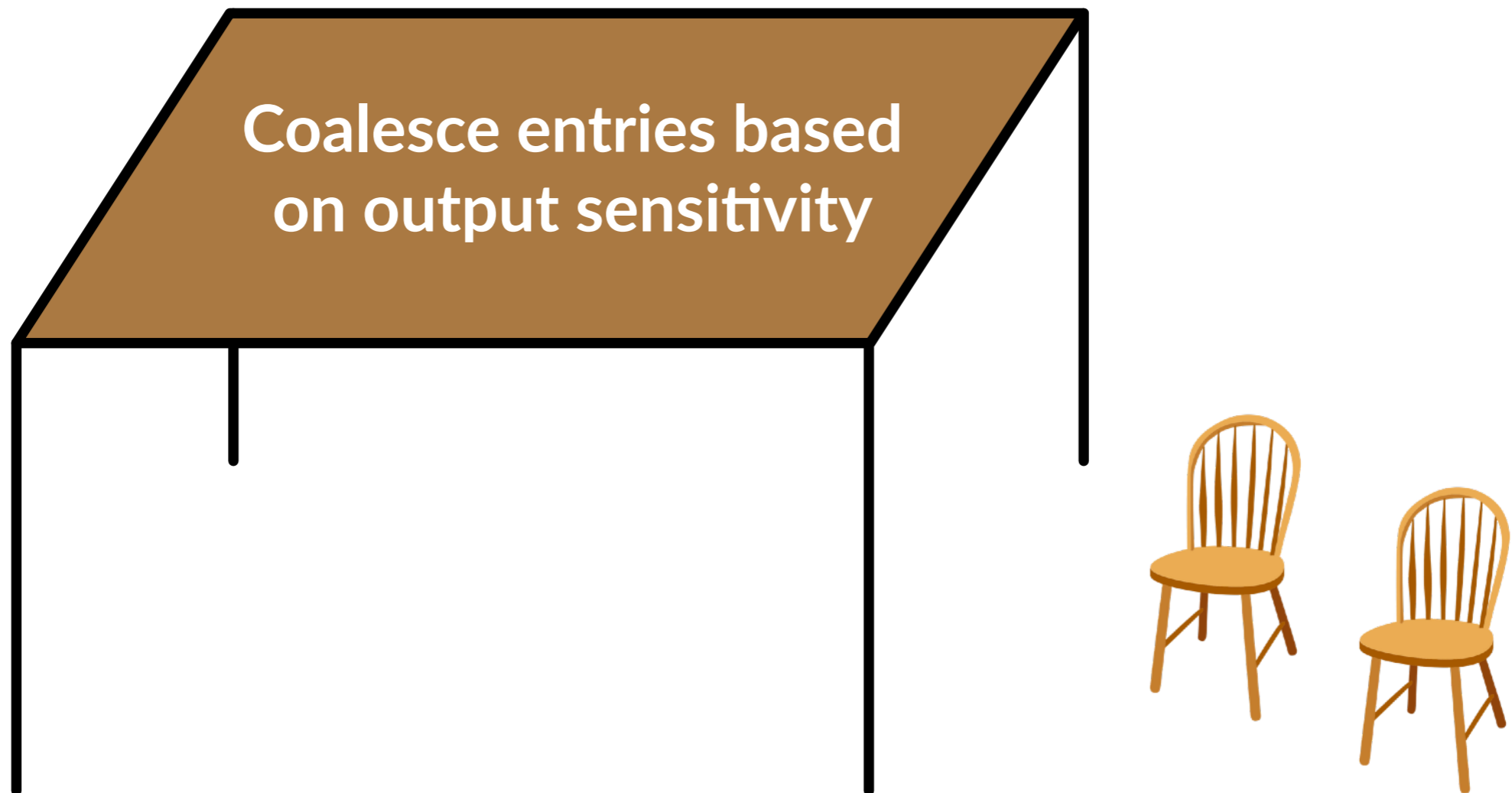
Coalesce entries based
on output sensitivity

Approximation can help...

There's prior work on this!

Fuzzy memoization for floating-point multimedia applications [TC '05]

Temporal approximate function memoization [IEEE Micro '18]

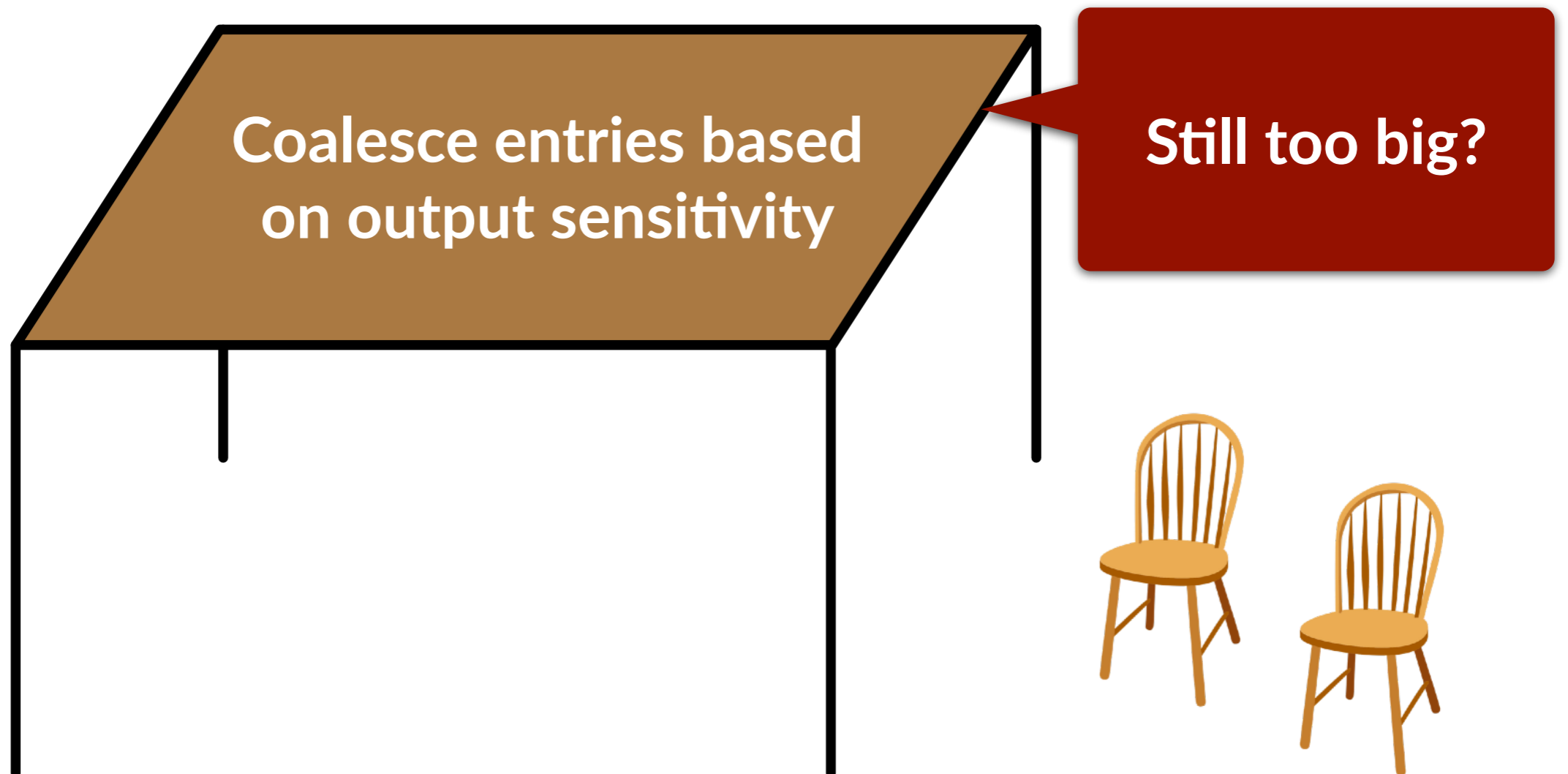


Approximation can help...

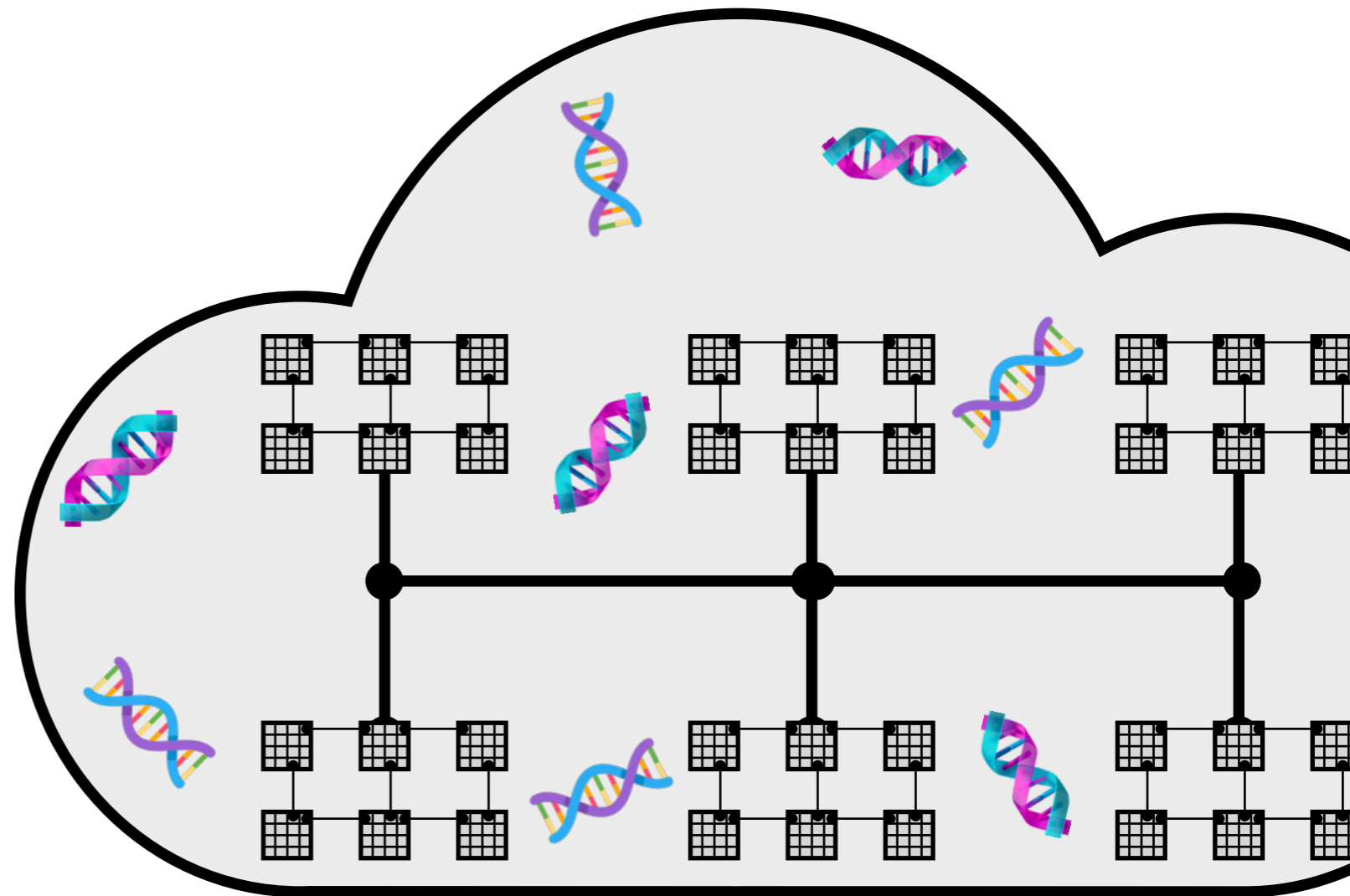
There's prior work on this!

Fuzzy memoization for floating-point multimedia applications [TC '05]

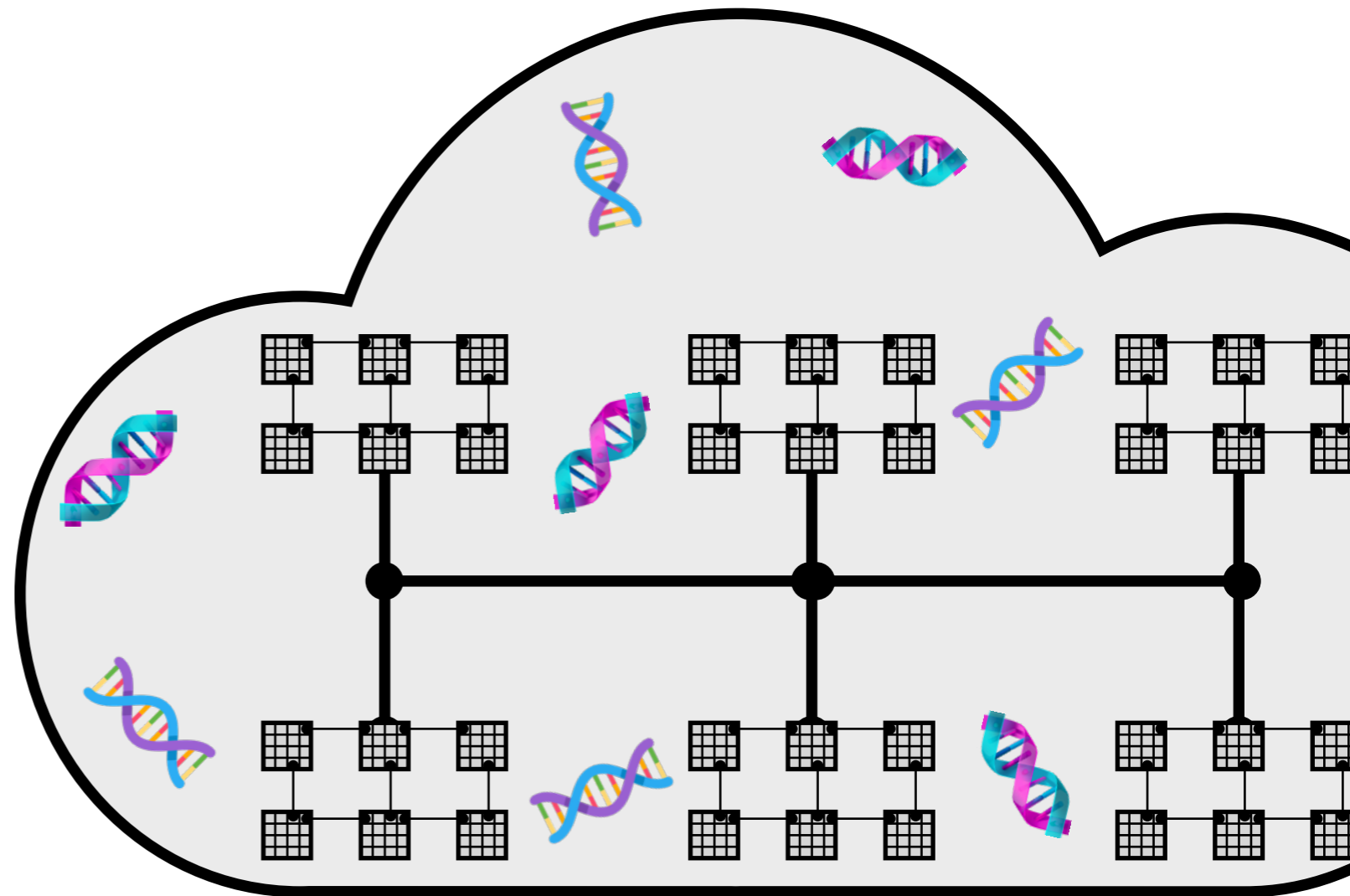
Temporal approximate function memoization [IEEE Micro '18]



Cloud-scale DNA storage to the rescue!



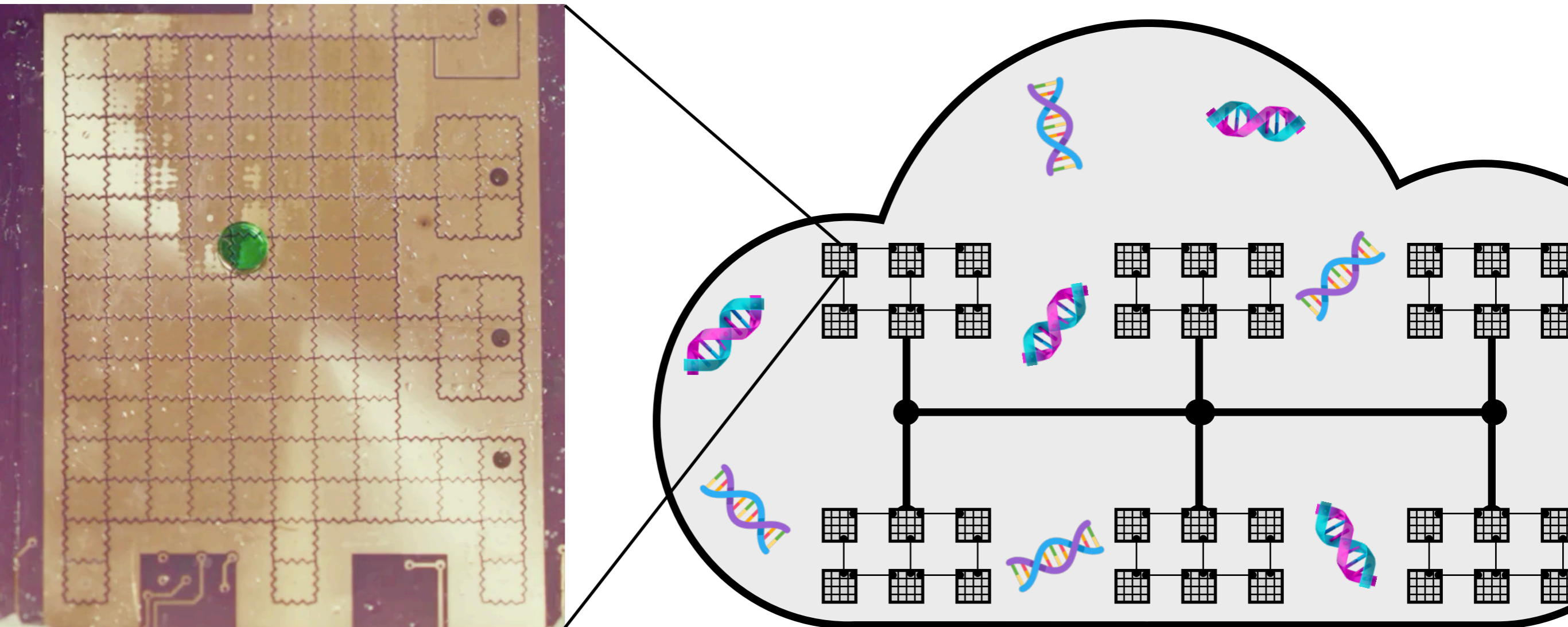
Cloud-scale DNA storage to the rescue!



Storing a Hellabyte would require *3 trillion* kg of SSDs,
but only 2500 kg of DNA!

Source: Max Willsey and Luis Ceze, *Mega-Microfluidics*, WACI 2019

Cloud-scale DNA storage to the rescue!



**Storing a Hellabyte would require 3 trillion kg of SSDs,
but only 2500 kg of DNA!**

Source: Max Willsey and Luis Ceze, *Mega-Microfluidics*, WACI 2019

Sharing is caring



Sharing is caring



Share $f(x)$
globally

Sharing is caring, *but is it safe?*



Share $f(x)$
globally

Sharing is caring, *but is it safe?*



Sharing is caring, *but is it safe?*



Sharing is caring, *but is it safe?*

Security and privacy implications

Attacker could brute force inputs to perform timing attacks on sensitive data

What happens to open-source code and cryptography?



A sampling of new research challenges



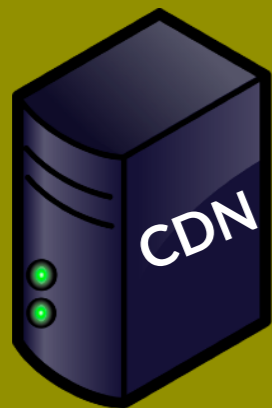
Cross-language
program equivalence



Security and privacy



Efficient networks



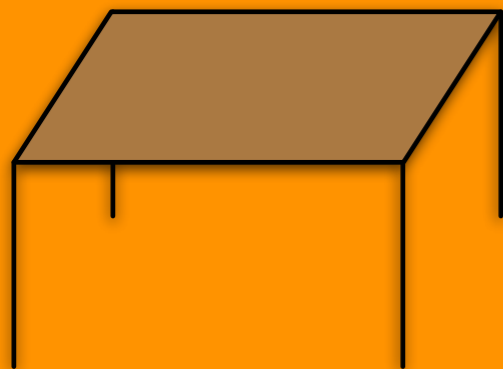
Global data delivery



High-density storage

#()

Large-scale hash
functions



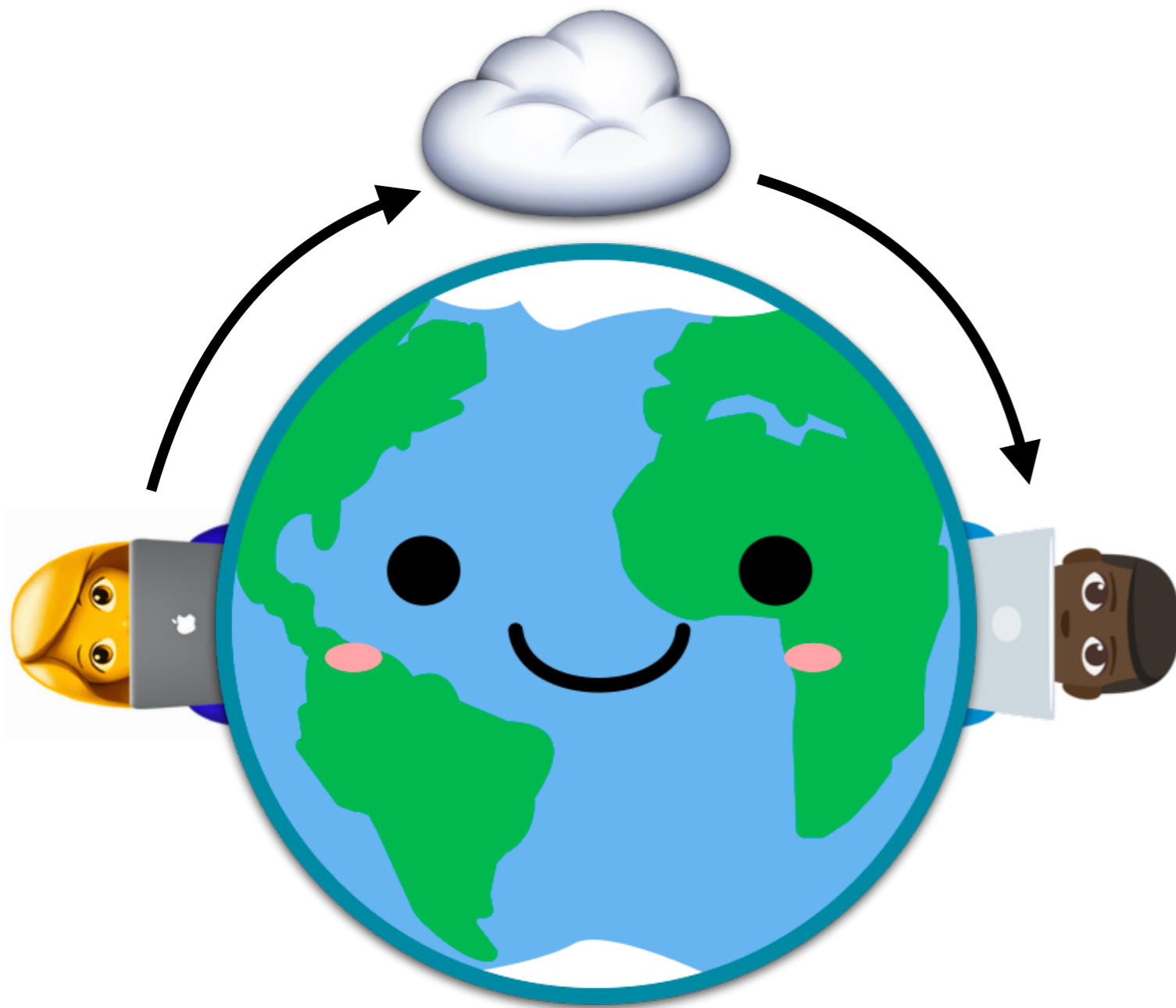
Fuzzy lookup tables



Accurate energy
models for functions

Your idea here!

Extreme Memoization Everything in a LUT!



Redundant
computation

Storing
results
locally

Extreme
memoization



Thanks!

Pratyush Patel — patelp1@cs.uw.edu