

The Datatypes Zoo

Gus Smith

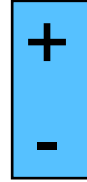
SAMPL Colloquium, 12/12/2019

By “datatypes”, we mean **numerical datatypes**: how the hardware represents and operates on real numbers.

IEEE 754 Floating Point

IEEE 754 Floating Point

Single precision (32 bit)



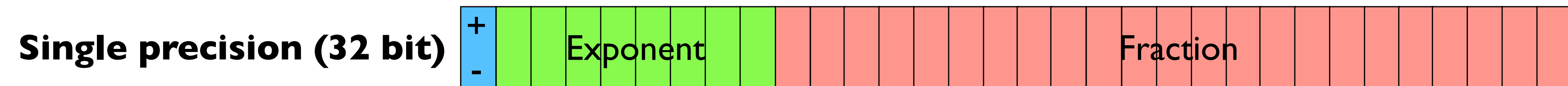
IEEE 754 Floating Point



IEEE 754 Floating Point

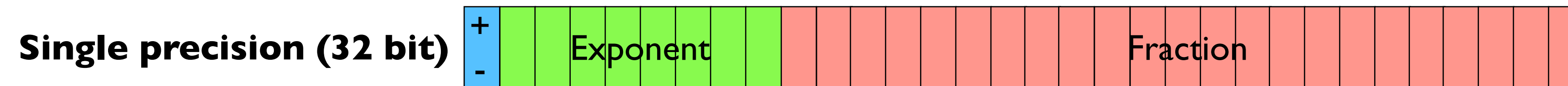


IEEE 754 Floating Point



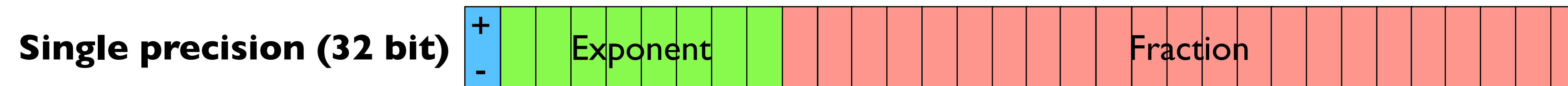
value \approx sign

IEEE 754 Floating Point



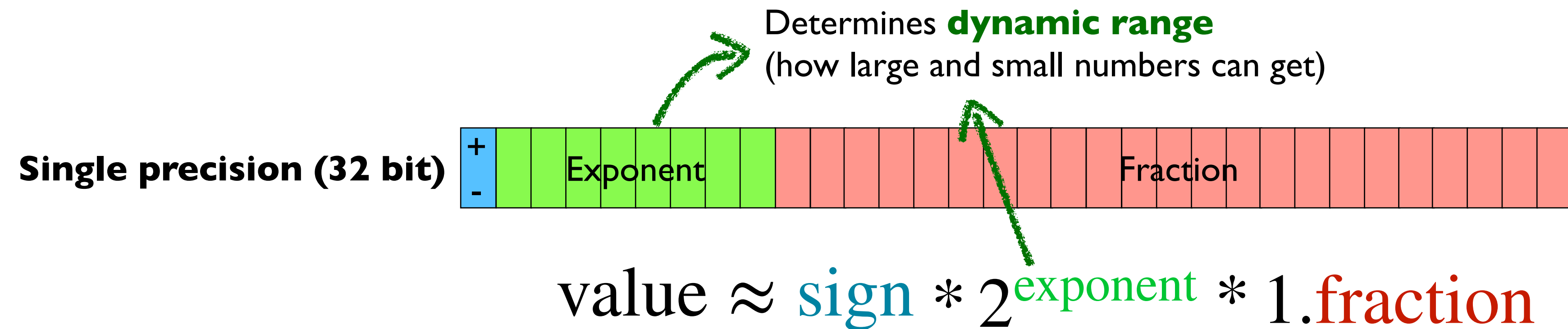
$$\text{value} \approx \text{sign} * 2^{\text{exponent}}$$

IEEE 754 Floating Point

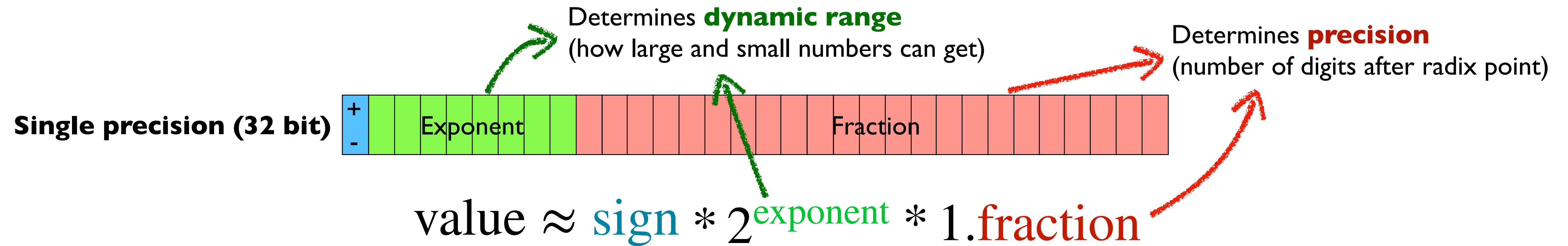


$$\text{value} \approx \text{sign} * 2^{\text{exponent}} * 1.\text{fraction}$$

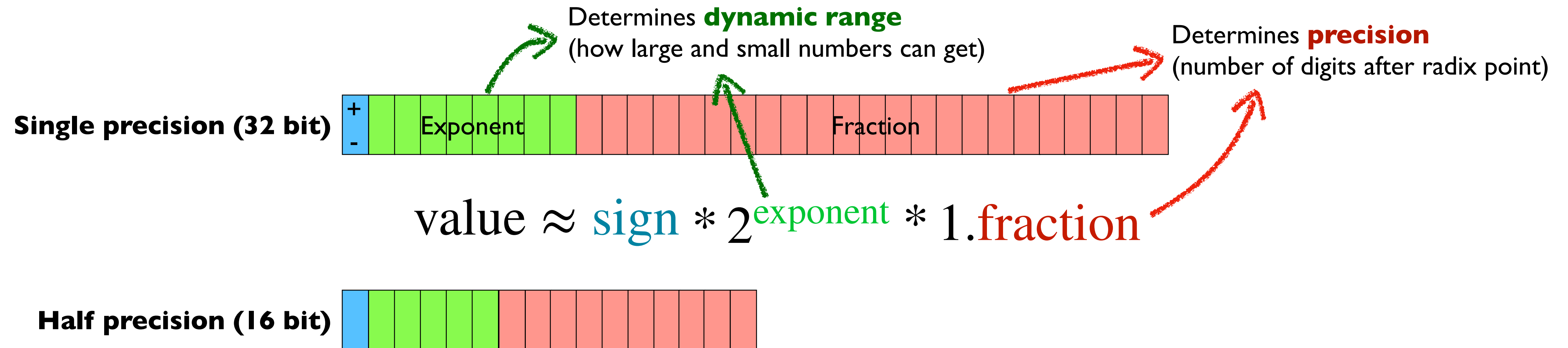
IEEE 754 Floating Point



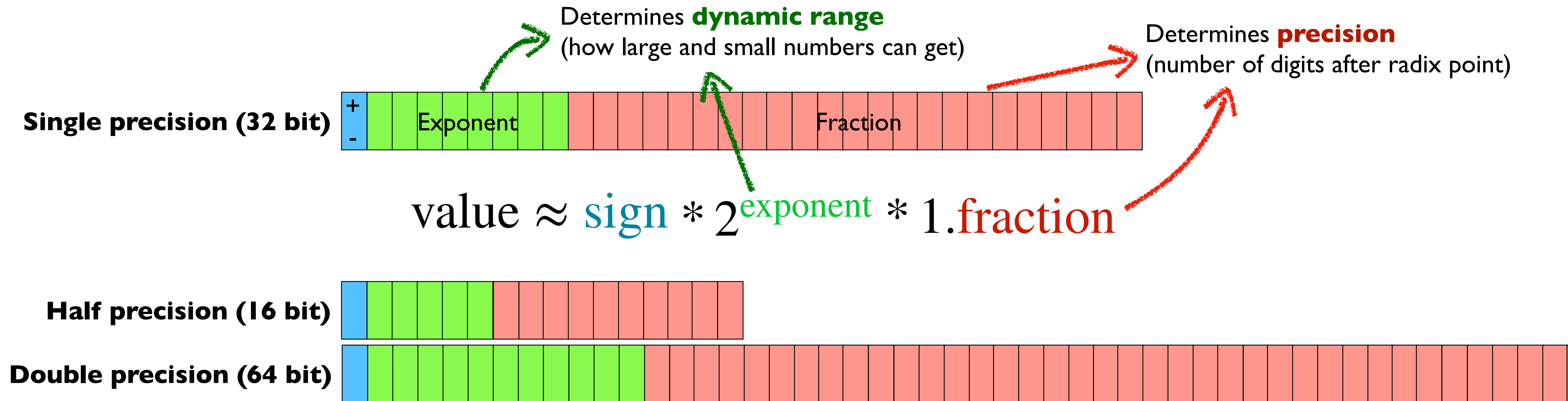
IEEE 754 Floating Point



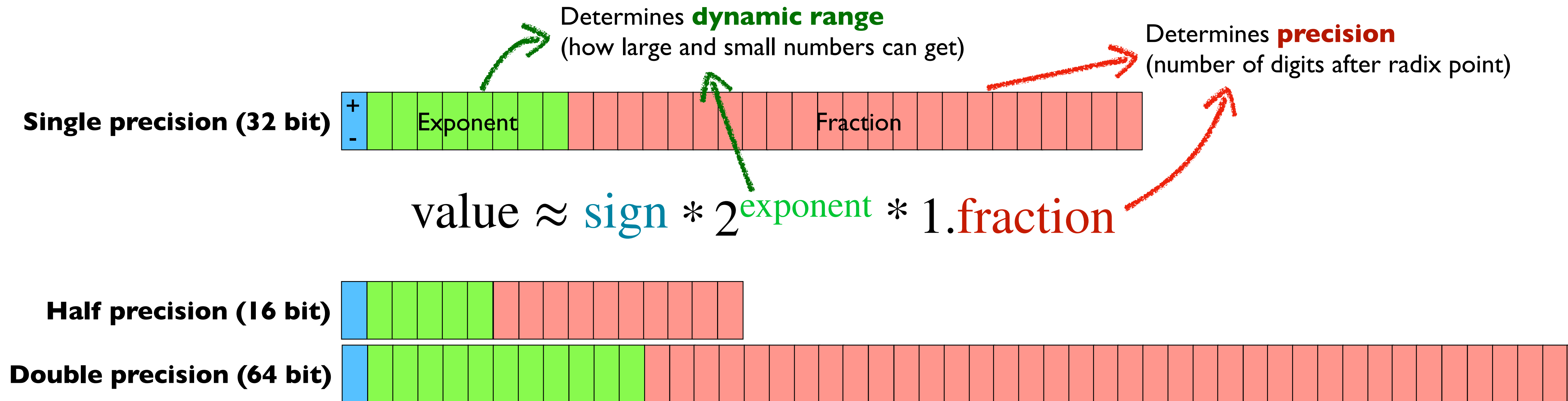
IEEE 754 Floating Point



IEEE 754 Floating Point

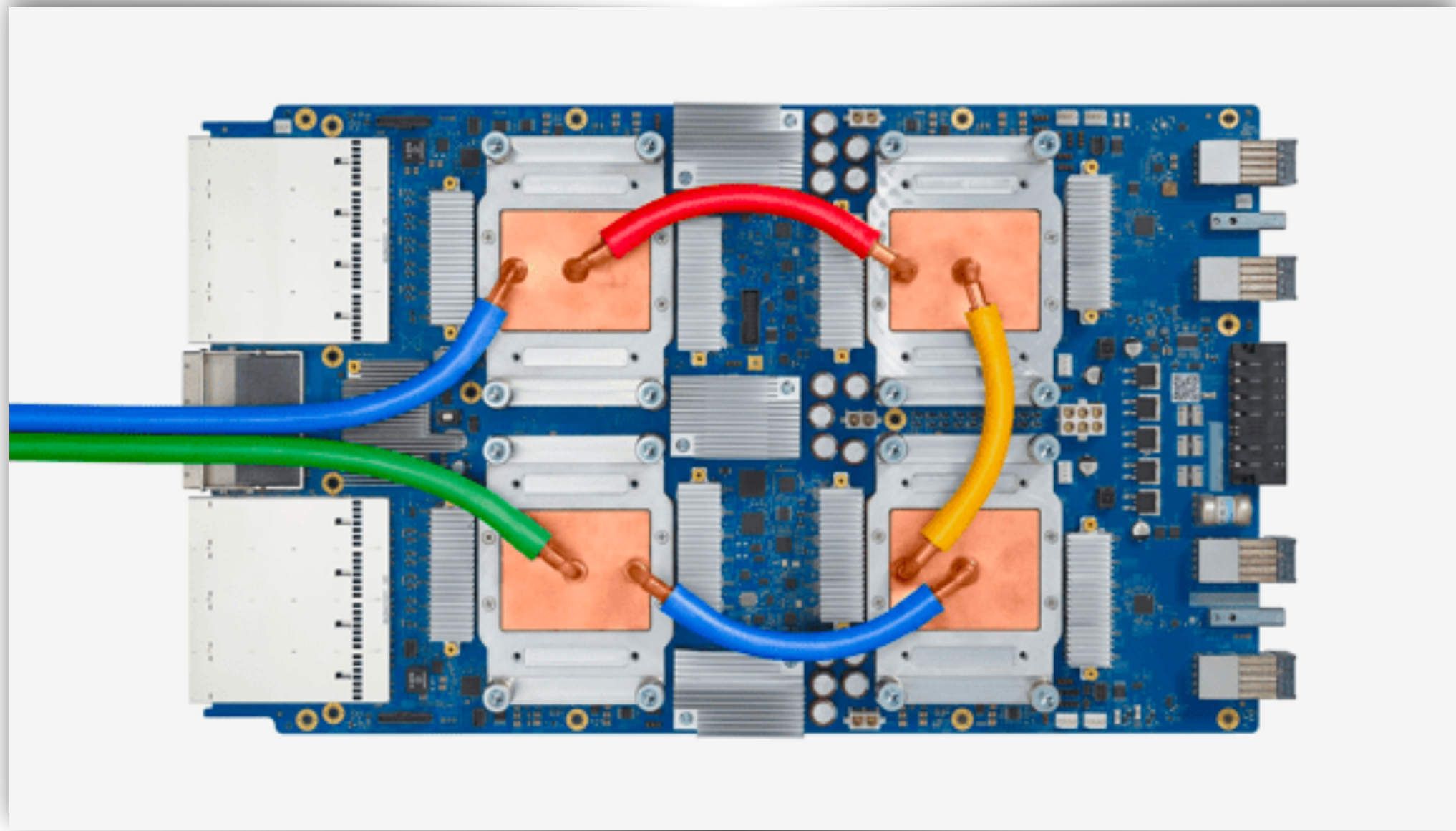


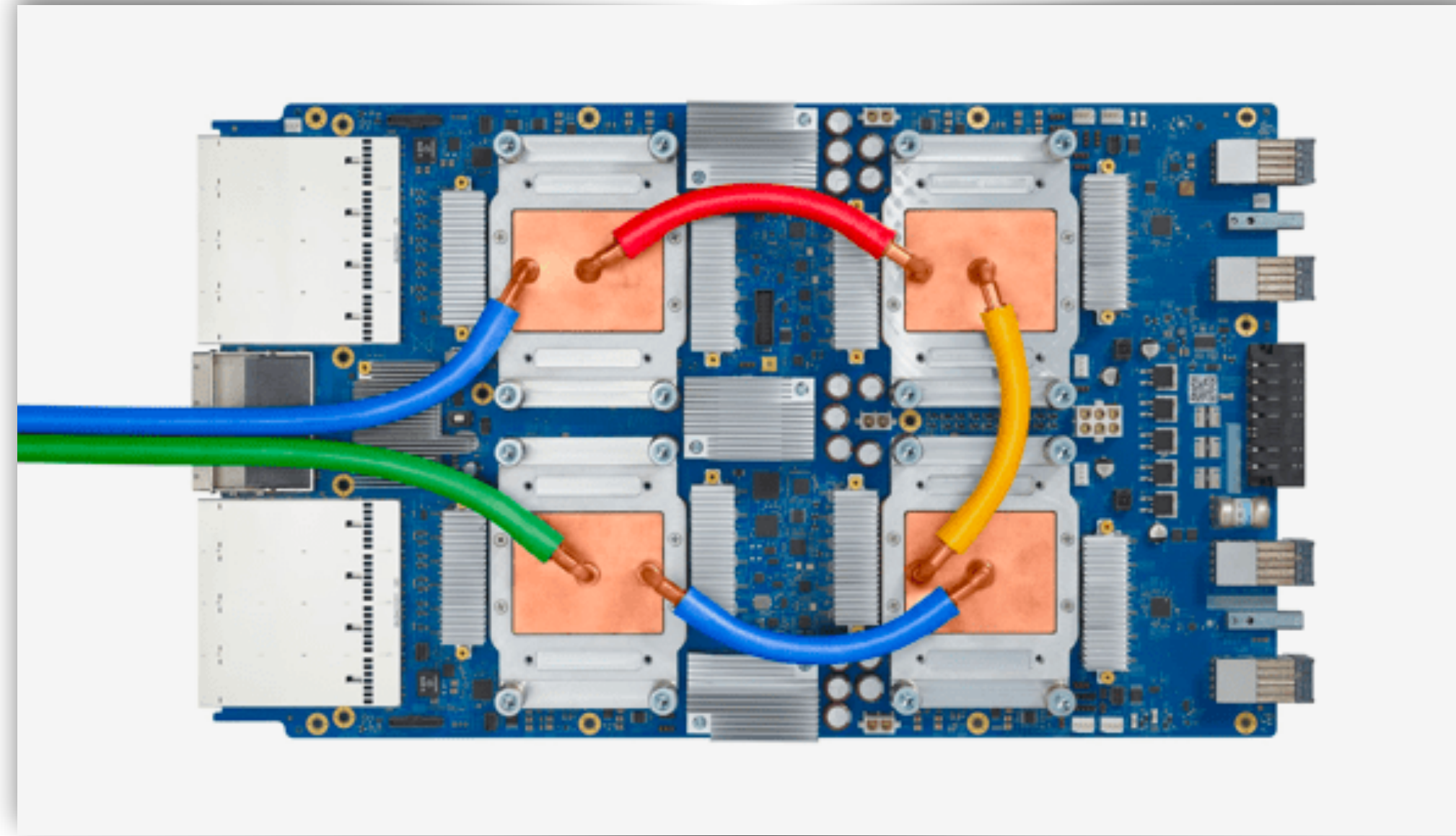
IEEE 754 Floating Point



Has remained an industry standard for more than thirty years!

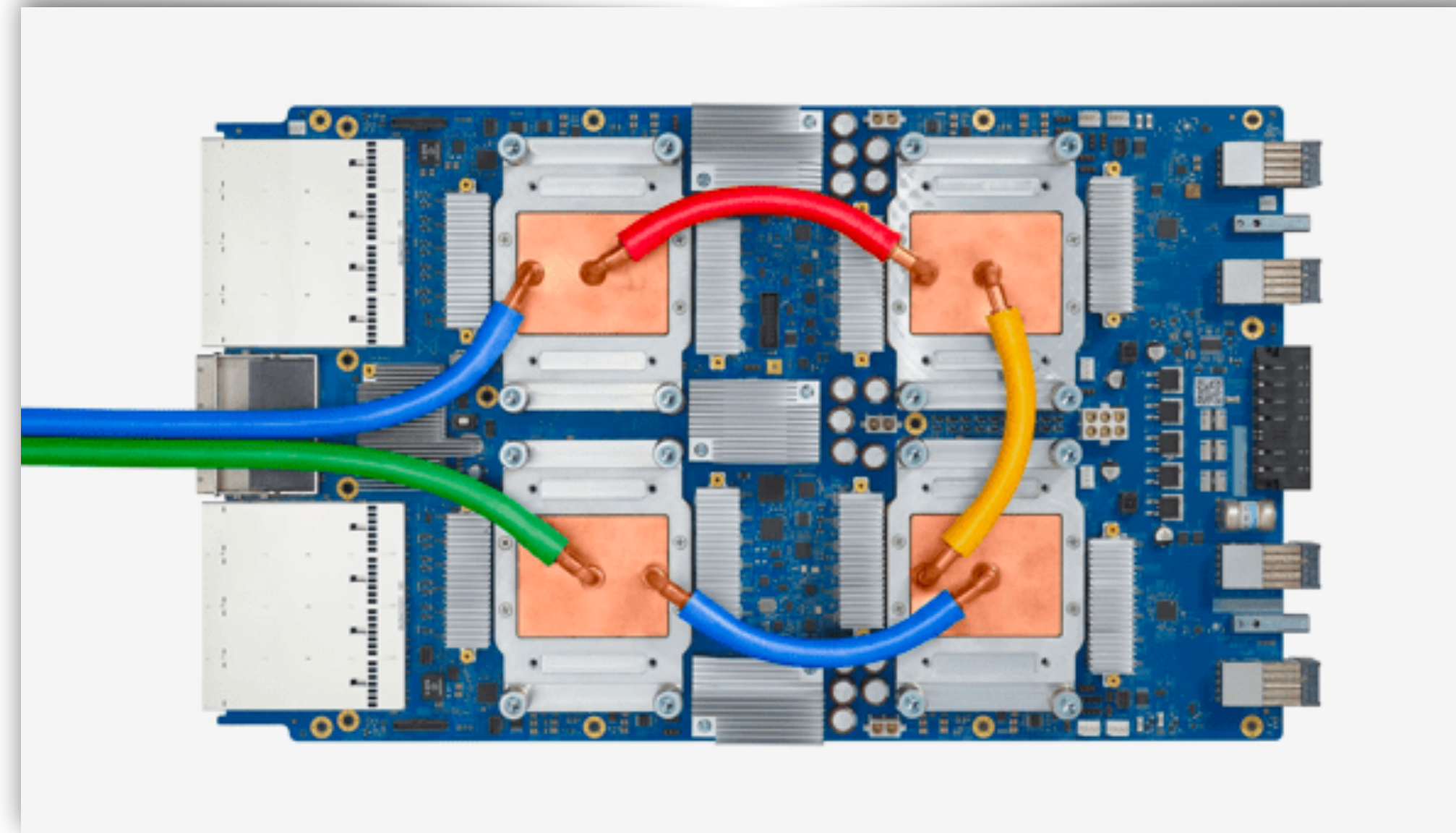






Should be fast and power-efficient

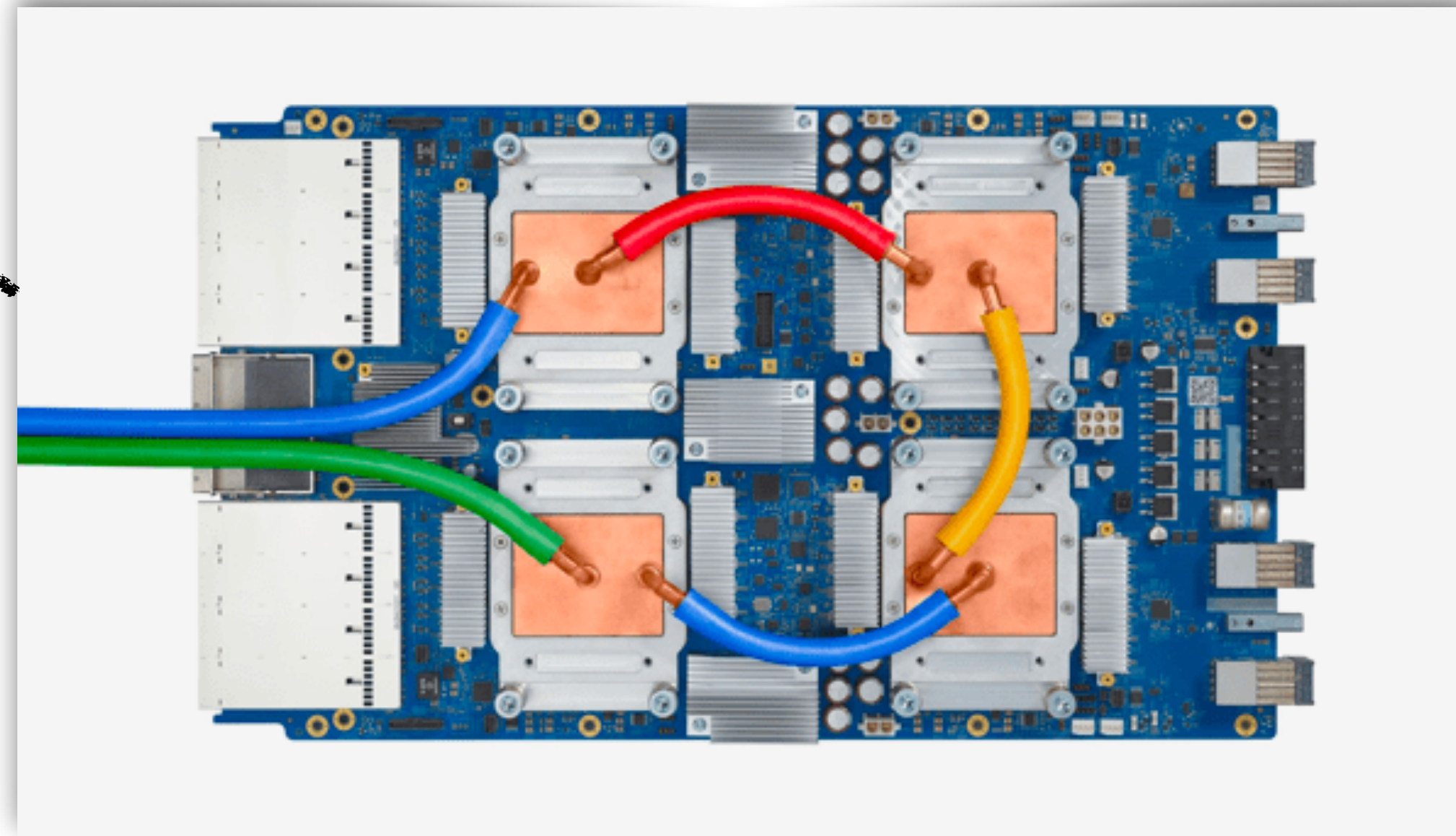




Should be fast and power-efficient

Needs small weights and activations
to maximize usage of chip area

Only needs to represent a specific range of values:
Weights and activations cluster (e.g. around $[-1, 1]$)

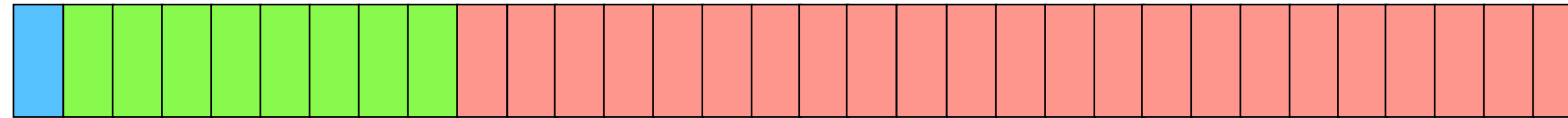


Should be fast and power-efficient

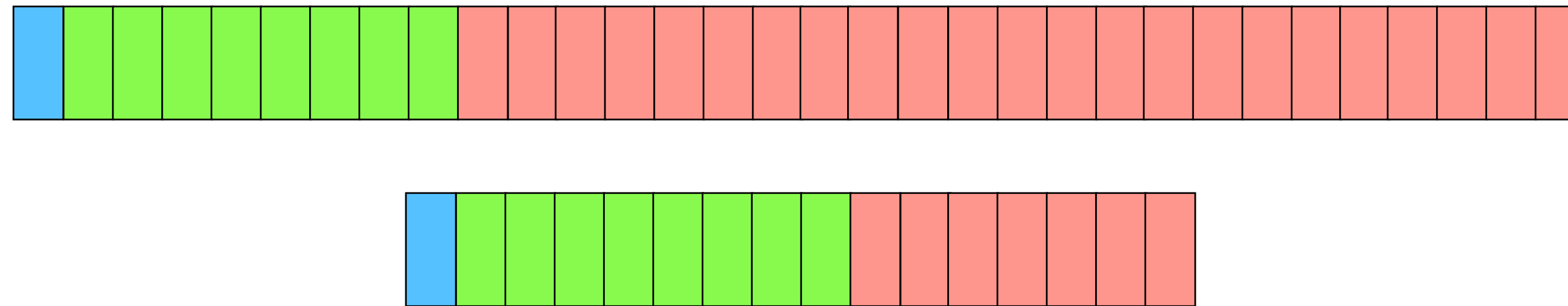
Needs small weights and activations
to maximize usage of chip area

The Datatypes Zoo

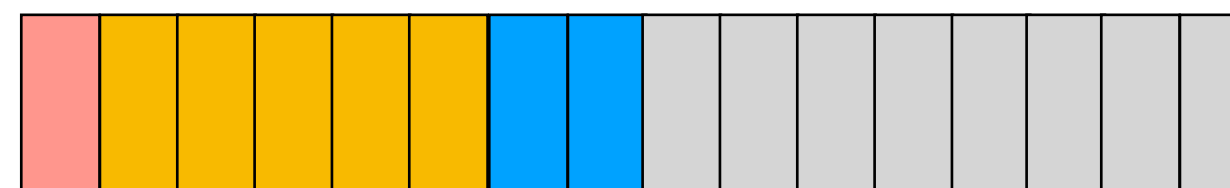
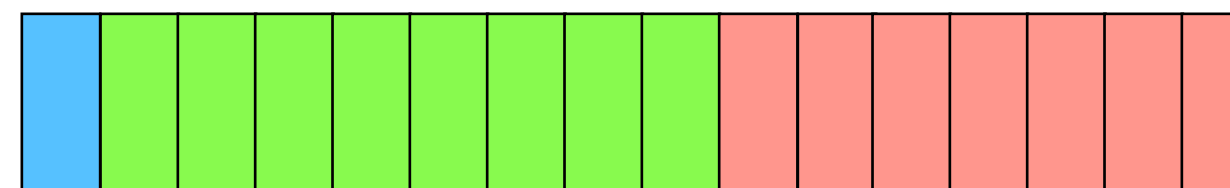
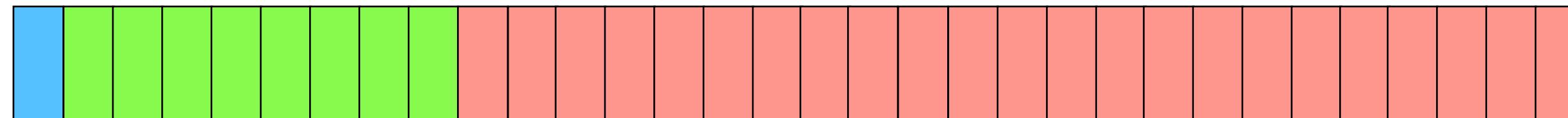
What do new datatypes look like?



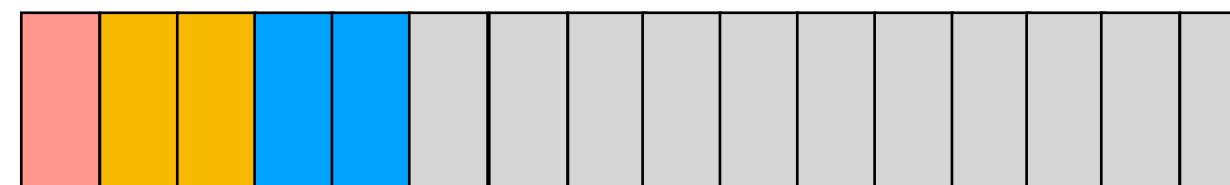
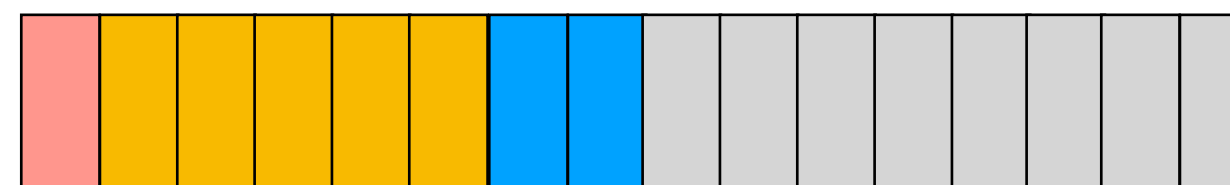
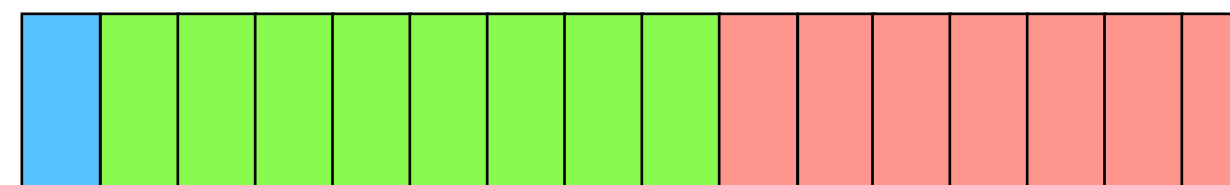
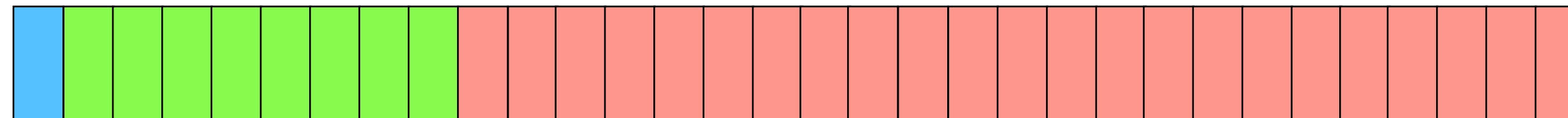
What do new datatypes look like?



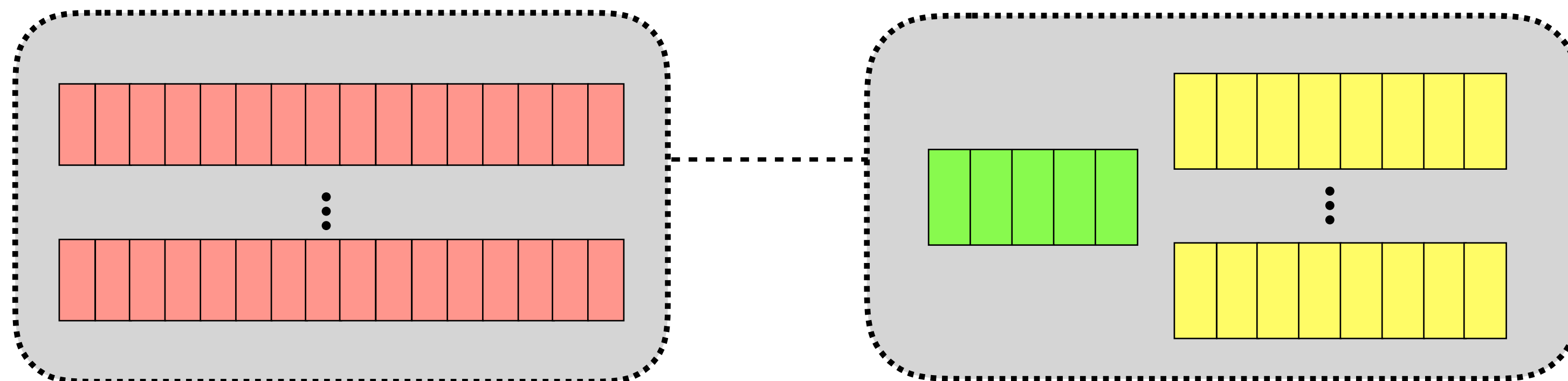
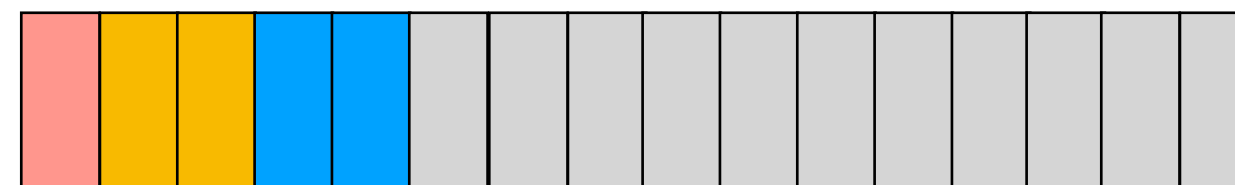
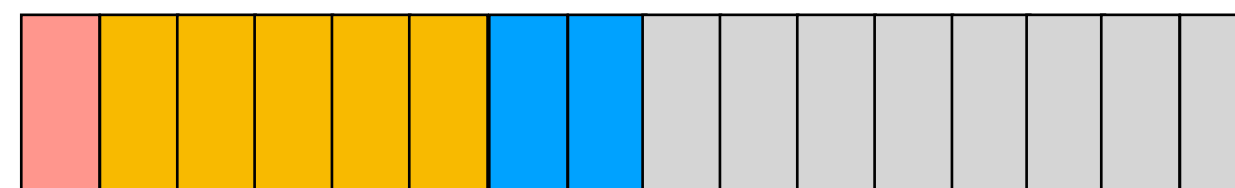
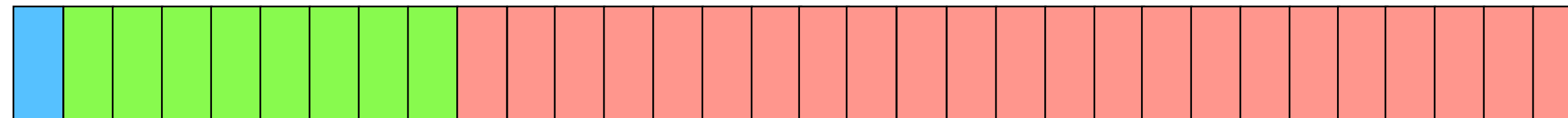
What do new datatypes look like?



What do new datatypes look like?

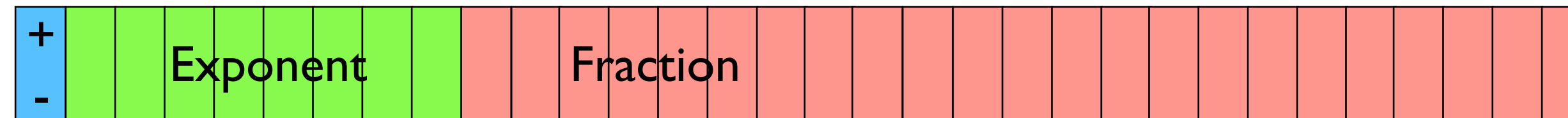


What do new datatypes look like?

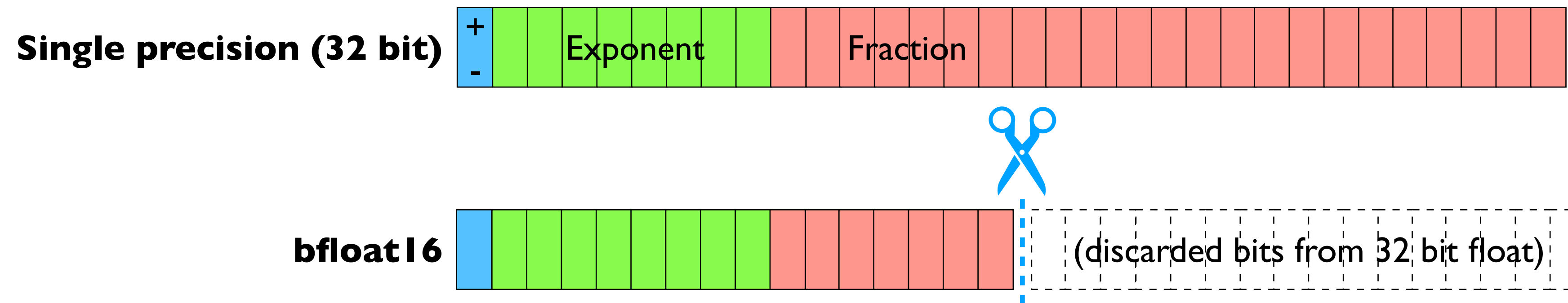


bfloat16

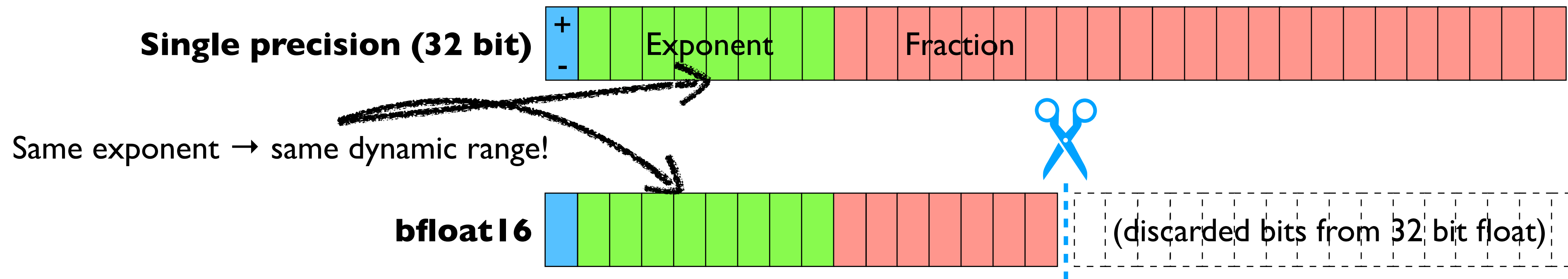
Single precision (32 bit)



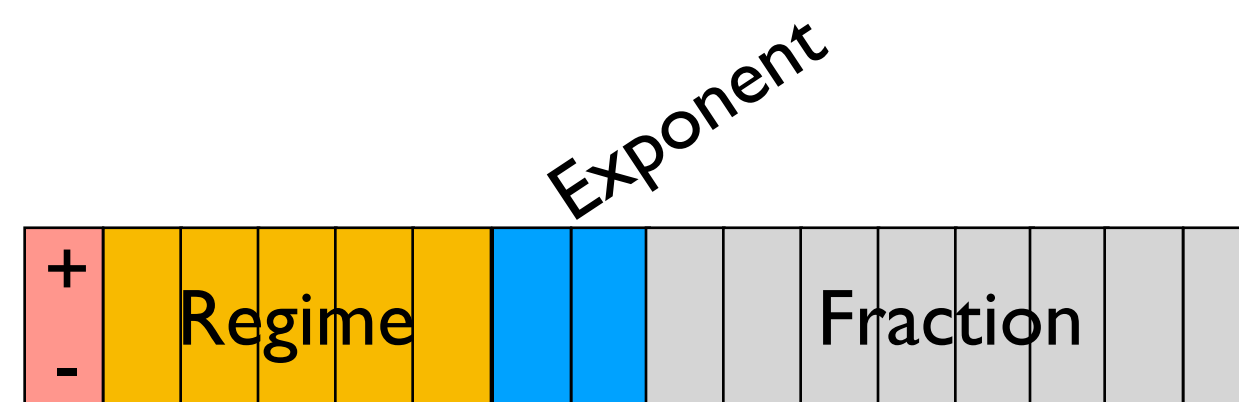
bfloat16



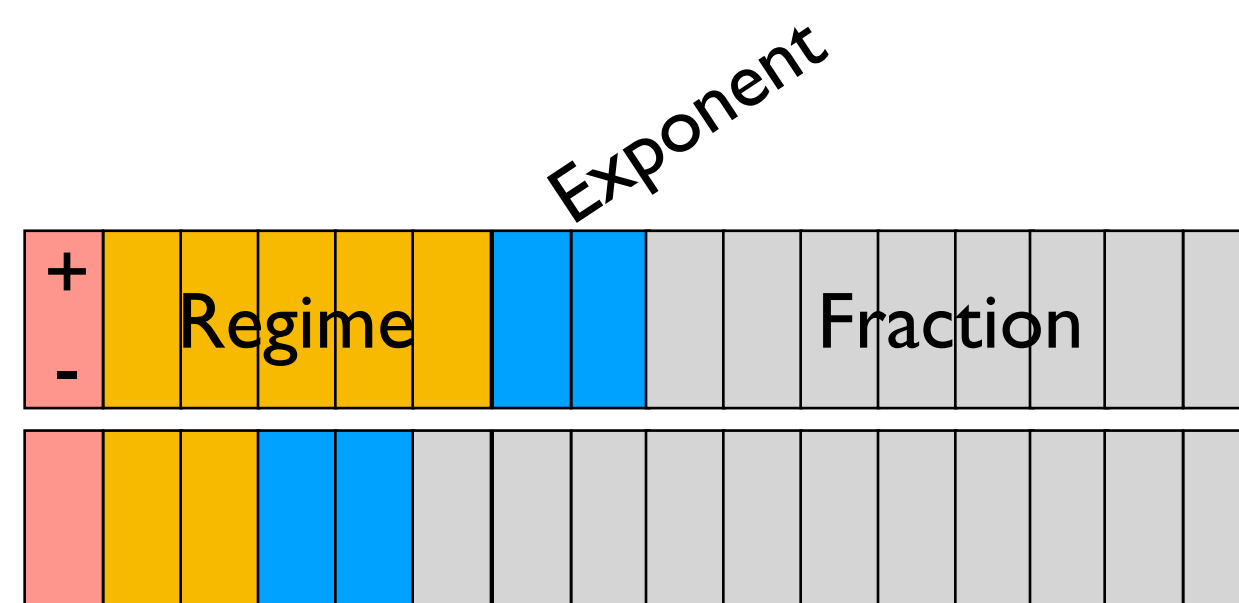
bfloat16



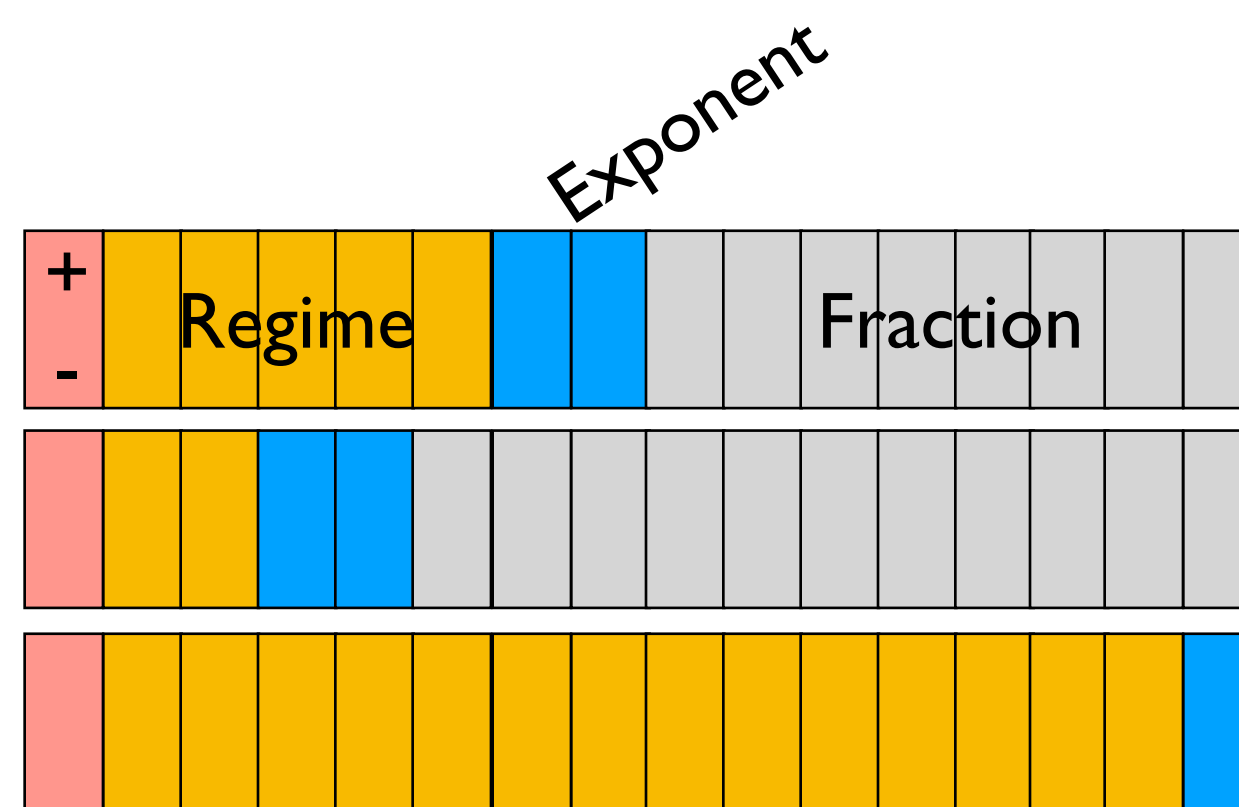
Posits



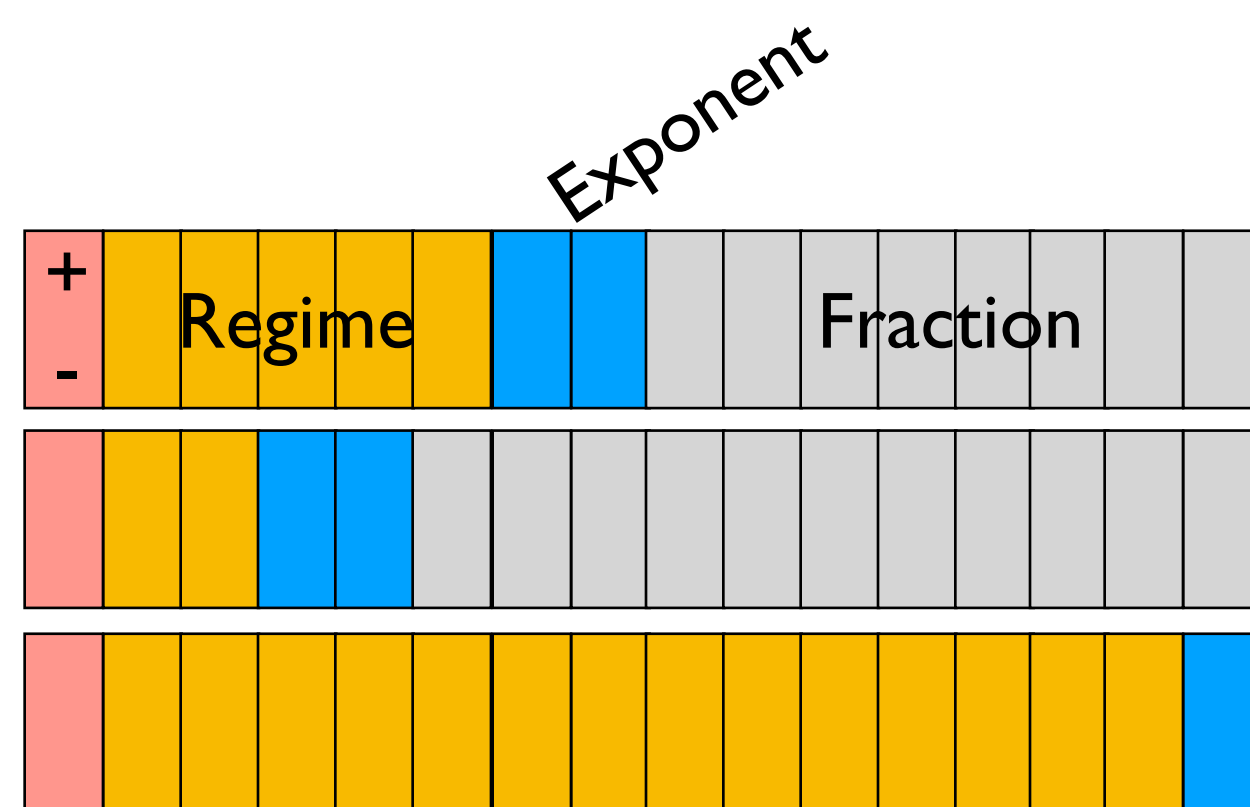
Posits



Posits

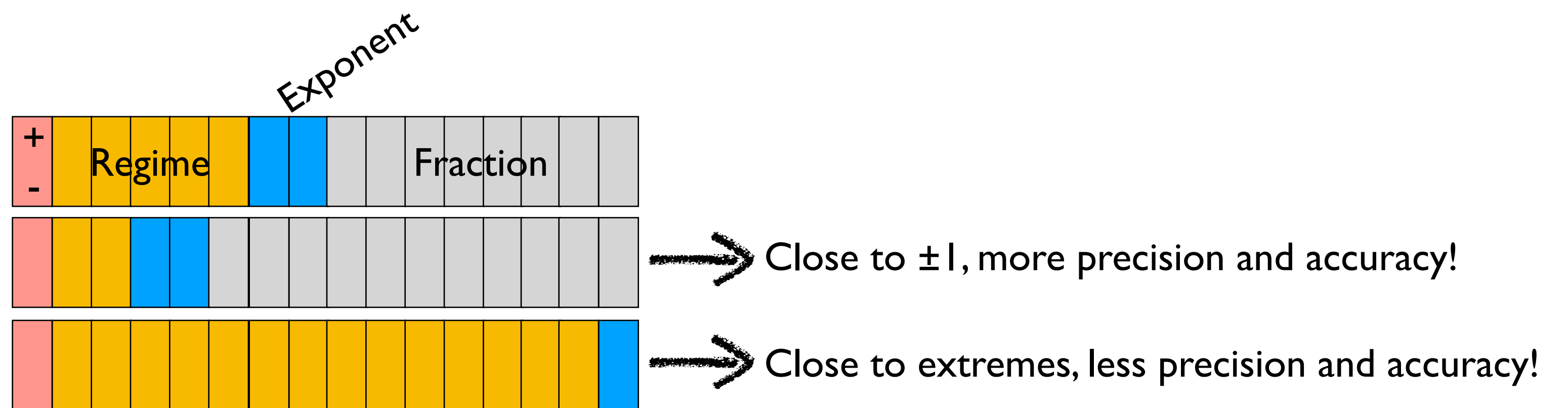


Posits

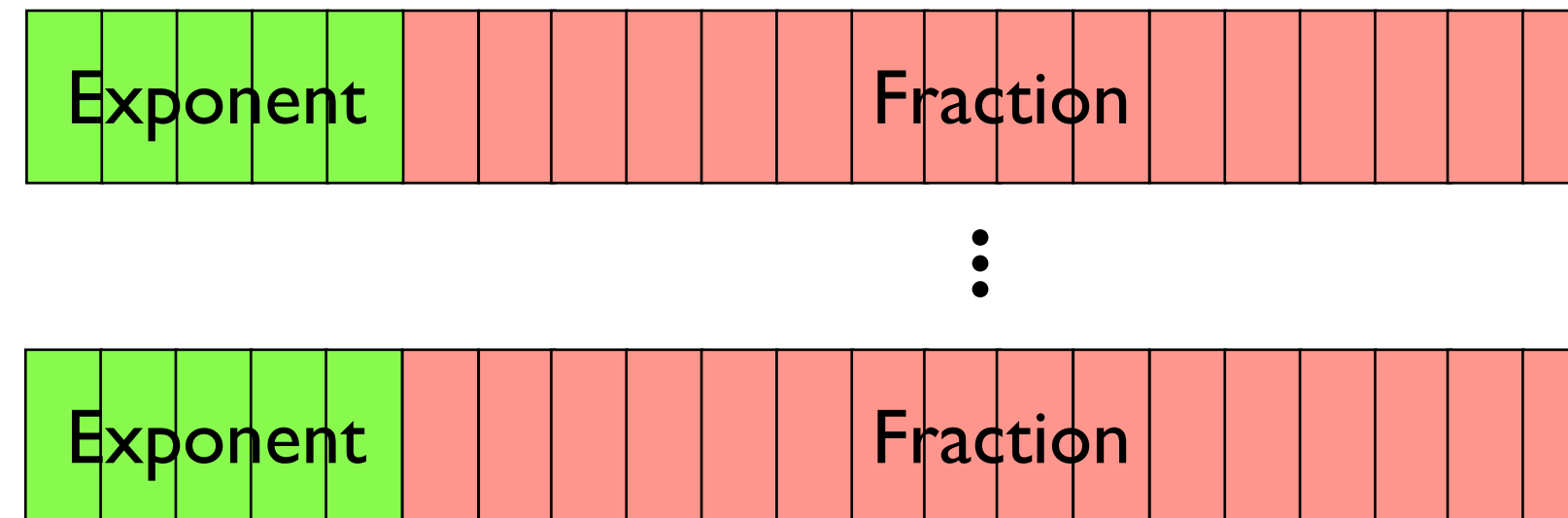


→ Close to ± 1 , more precision and accuracy!

Posits

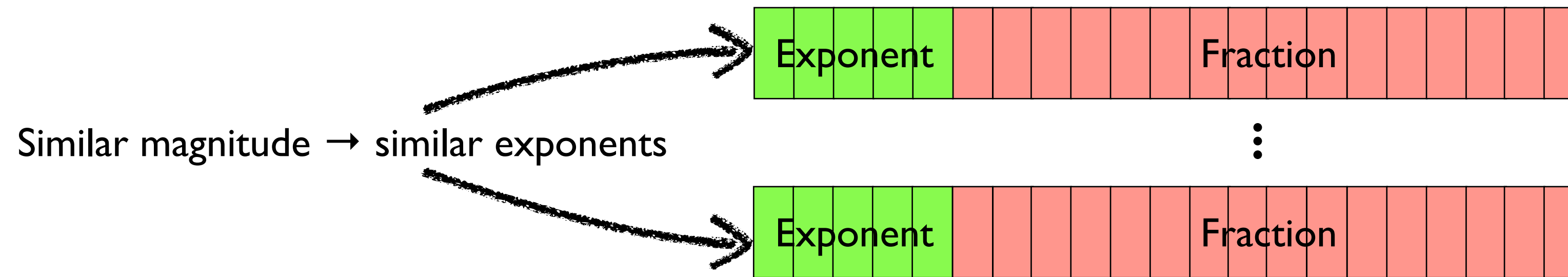


Flexpoint

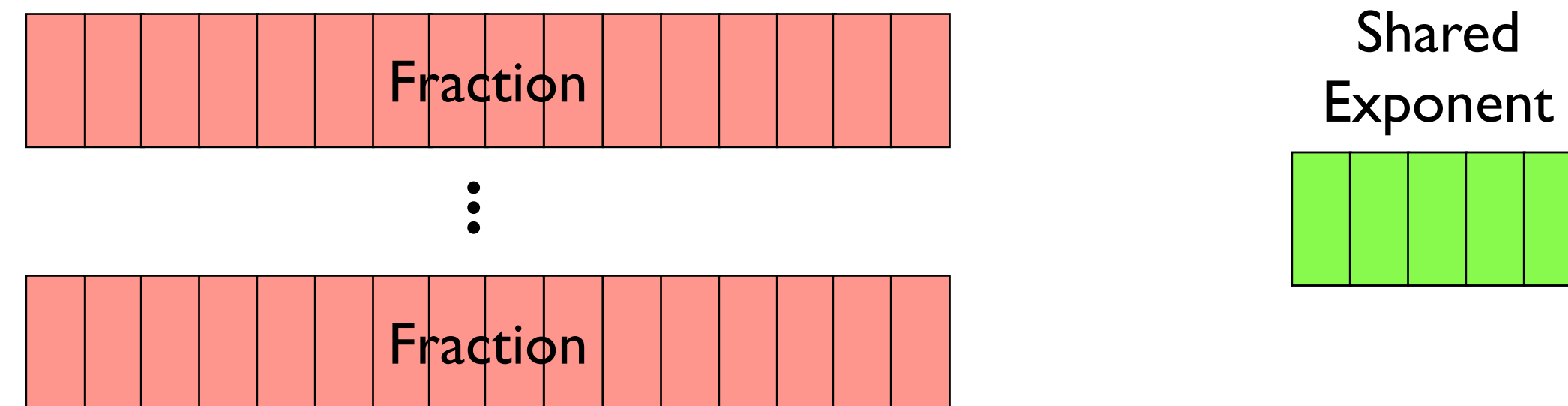


See Köster et. al., [“An Adaptive Numerical Format for Efficient Training of Deep Neural Networks”](#)

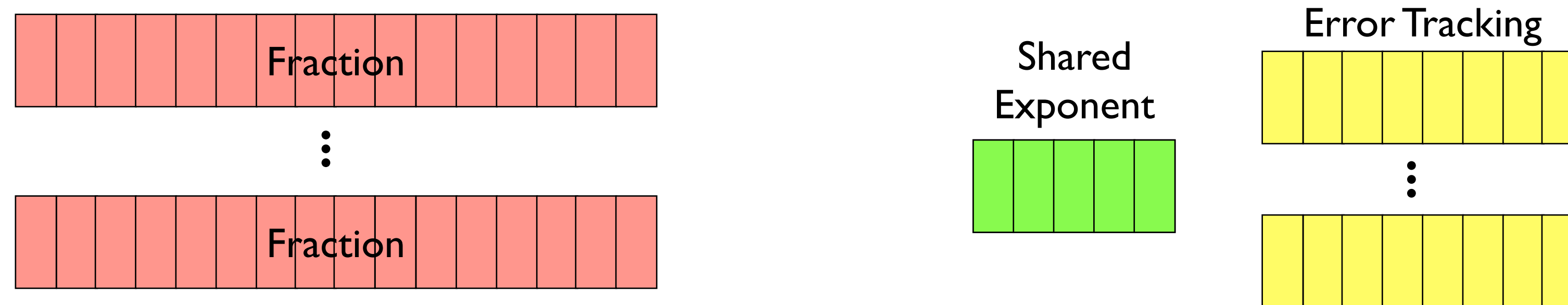
Flexpoint



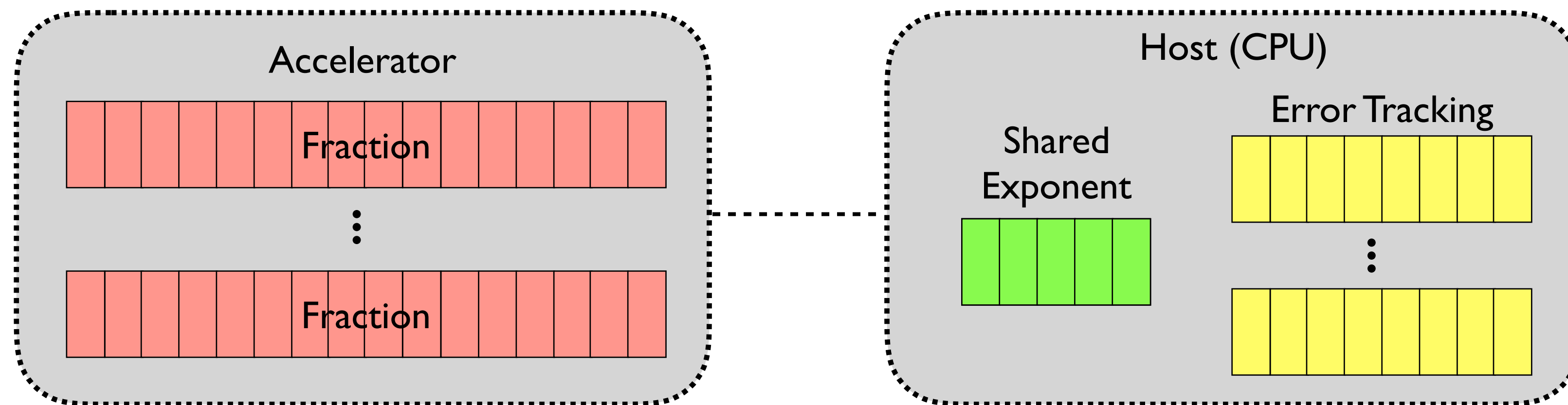
Flexpoint



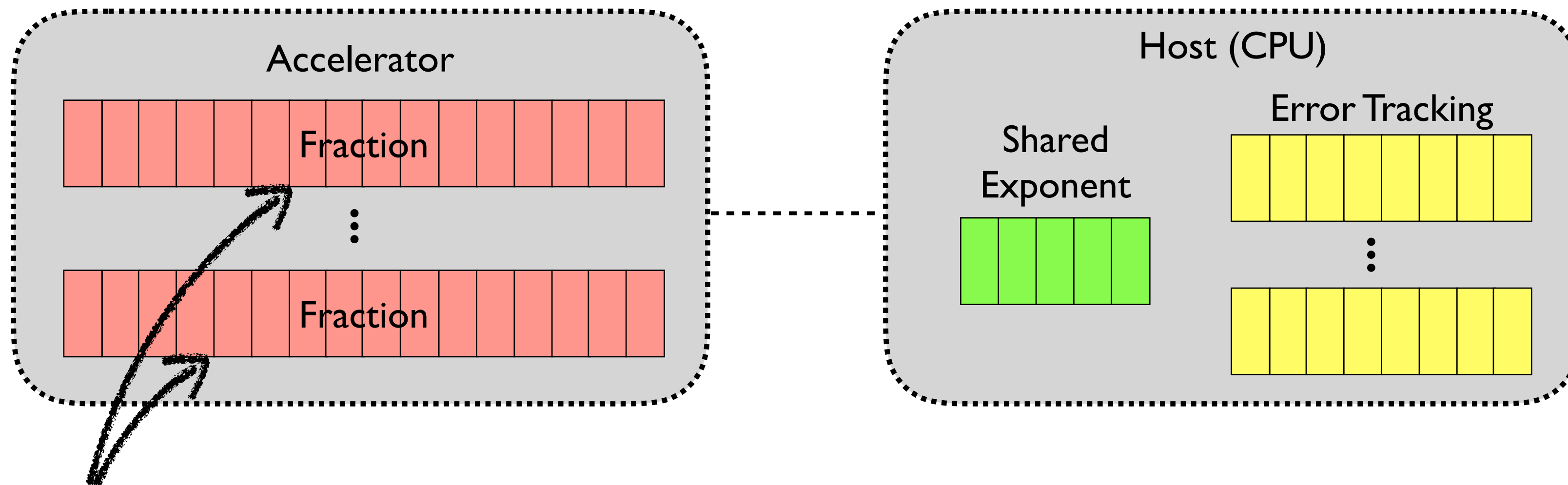
Flexpoint



Flexpoint



Flexpoint



Just integers! We can use integer hardware!

In conclusion,

In conclusion,

- Representing real numbers in hardware is an ongoing challenge

In conclusion,

- Representing real numbers in hardware is an ongoing challenge
- There are many interesting solutions out there, beyond IEEE floats!

Thank you!