

Gables: A Roofline Model for Mobile SoCs

Mark D. Hill, Wisconsin & Former Google Intern
Vijay Janapa Reddi, Harvard & Former Google Intern

One Roofline



HPCA
Feb 2019

ASPLOS
Tutorial

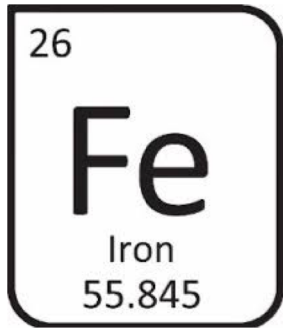
N Rooflines



Gables: A Roofline Model for Mobile SoCs

Mark D. Hill, Wisconsin & Former Google Intern
Vijay Janapa Reddi, Harvard & Former Google Intern

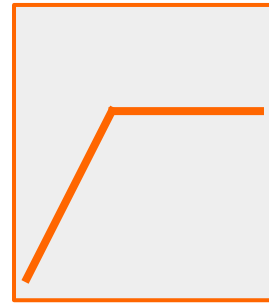
Models give insight & first answer



**CPU &
Iron Law**



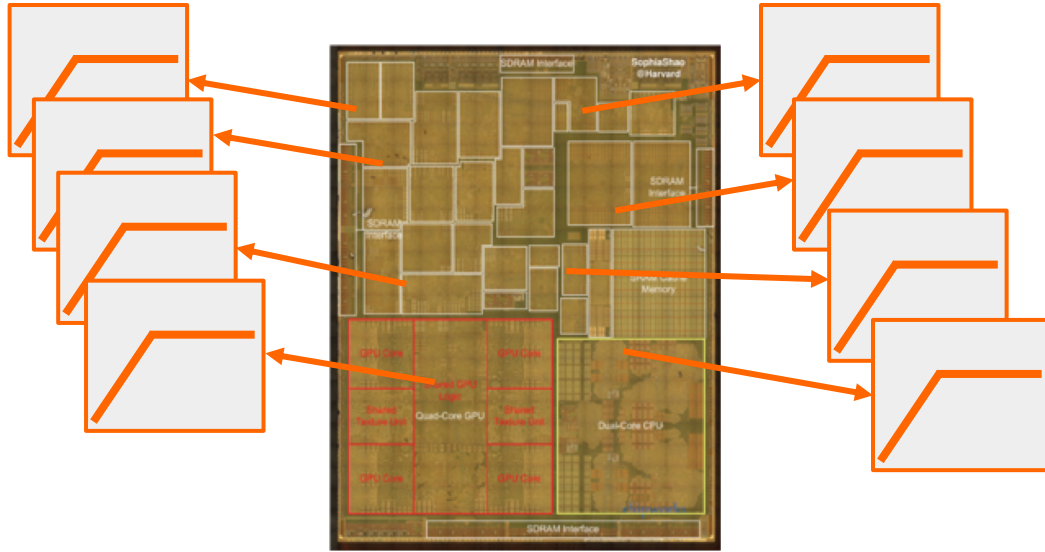
**Multiprocessor &
Amdahl's Law**



**Multicore &
Roofline**

**Mobile
System on a
Chip (SoC)
&
What???**

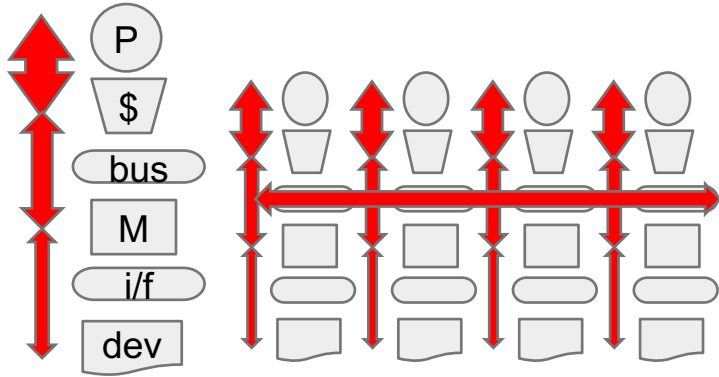
Mobile System on Chip (SoC) & Gables



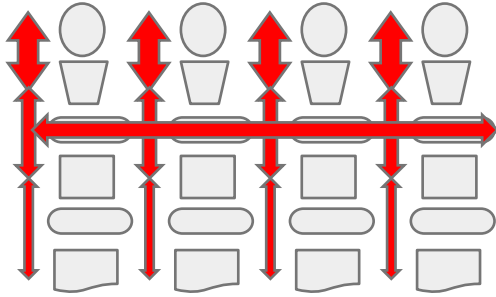
Dozens of Accelerators. Which?
Provision “Goldilocks” Accelerators?
Whither workload communication?

Gables @
~5:55pm
today!

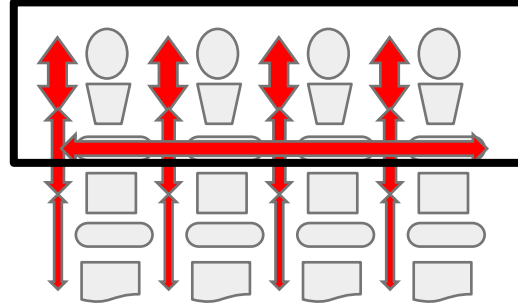
Computer Architecture & Models



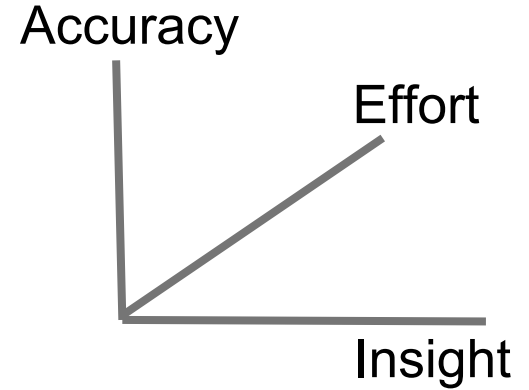
**CPU &
Iron Law**



**Multiprocessor &
Amdahl's Law**



**Multicore &
Roofline**

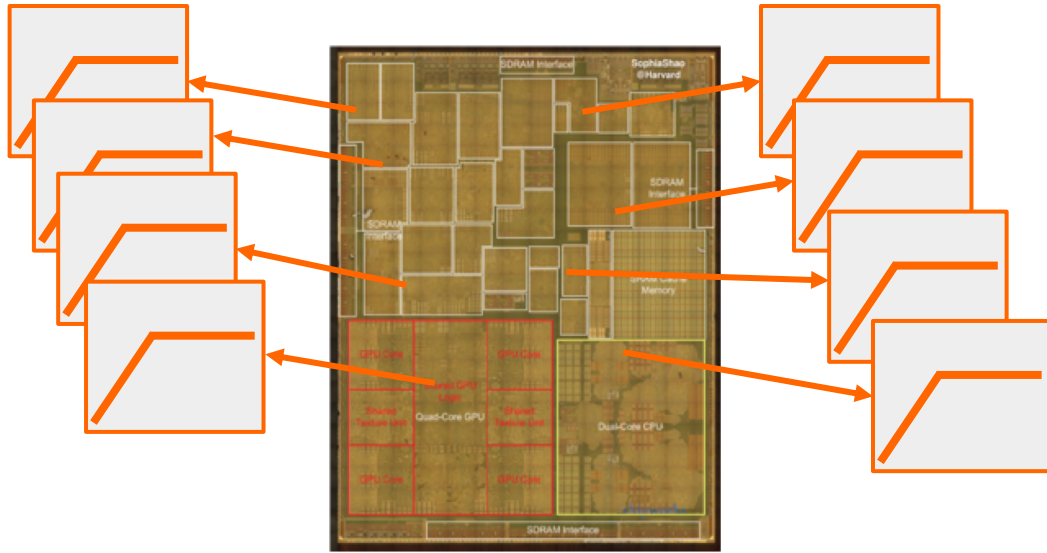


Models v.v. Simulation

- **More insight**
- **Less effort**
- But less accuracy

Models give first answer, not final answer

Mobile System on Chip (SoC) & Gables



1. Include Accelerator IP[i]?
2. IP[i] over-provisioned?
3. IP[i] over-communicates?

Gables provides first answer!

Outline

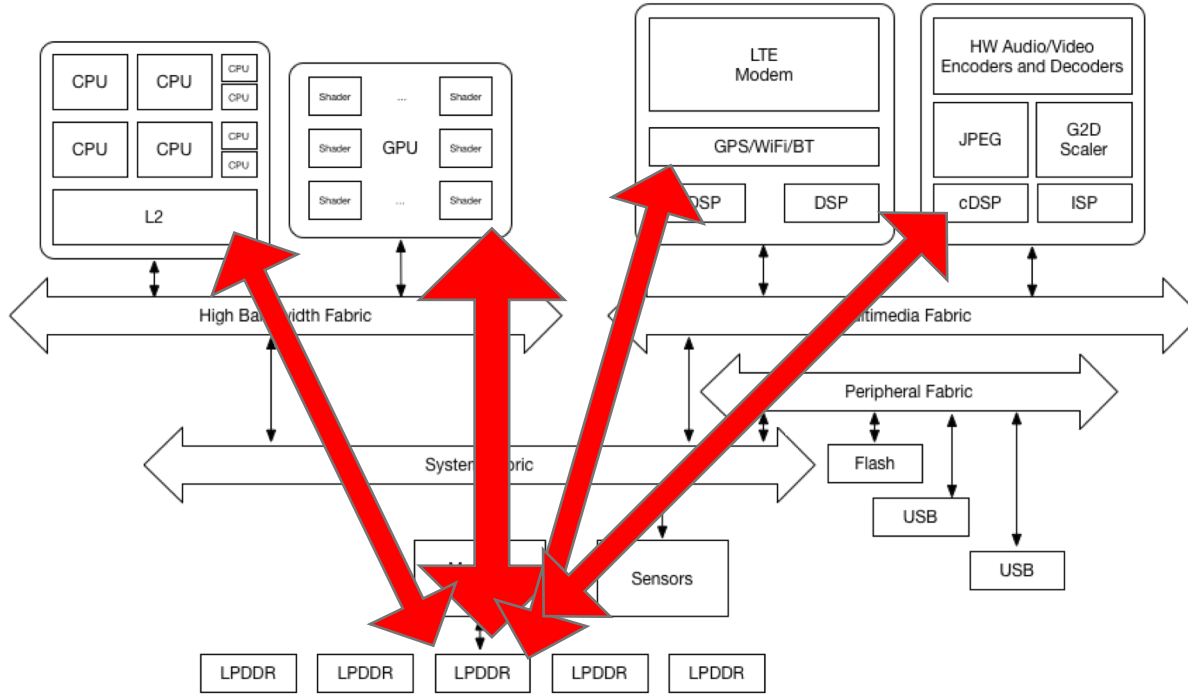
Motivation: Mobile SoCs, Usecases, & IP Blocks

Gables: A Roofline Model for Mobile SoCs

Gables 2-IP Example

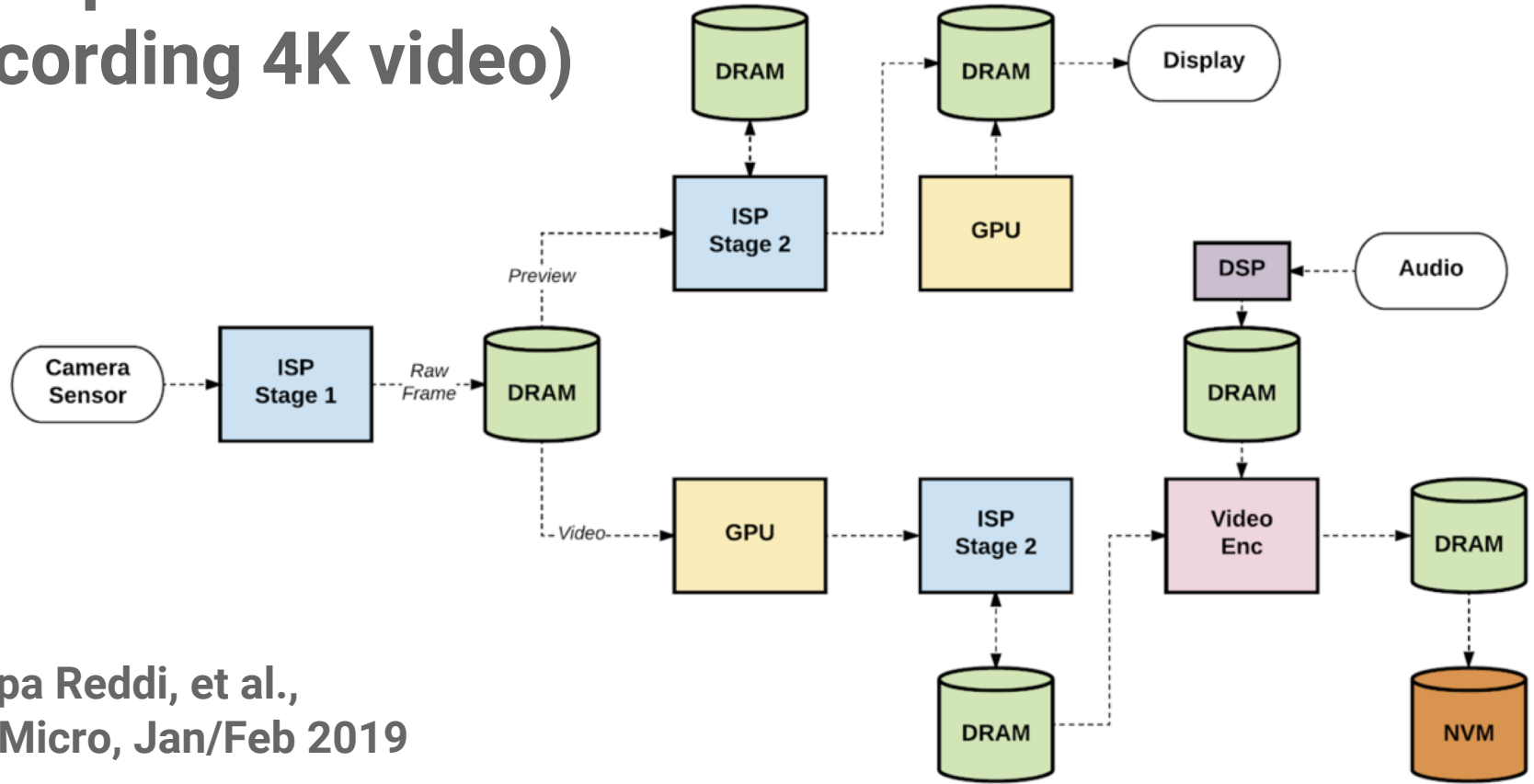
Wrap Up

Mobile SoC HW



Many IP blocks; Many flows, Many degrees of freedom

Example Usecase (recording 4K video)



Janapa Reddi, et al.,
IEEE Micro, Jan/Feb 2019

Mobile SoCs Run Usecases

	AP	Display	G2DS	GPU	ISP	JPEG	IPU	VDEC	VENC	DSP
HDR+	X	X		X	X	X	X			
Videocapture	X	X		X	X				X	
VideocaptureHDR	X	X		X	X				X	
VideoplaybackUI	X	X	X	X				X		
Google Lens	X	X	X	X						X

Must run each usecase sufficiently fast -- no need faster

Must run all usecases – average irrelevant

A usecase uses IPs concurrently – more than serially

A. Commercial SoCs Hard To Design

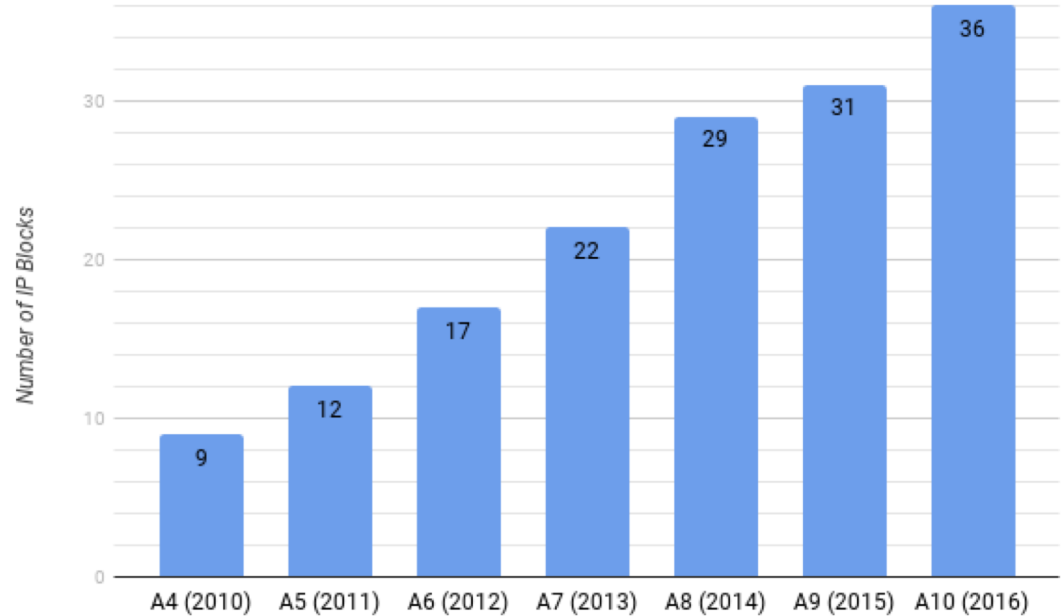
Envision usecases
(2-3 years ahead)

Select IPs

Size IPs

Design Uncore

Cycle-level simulation
later, but....



What about early before SW written?

Outline

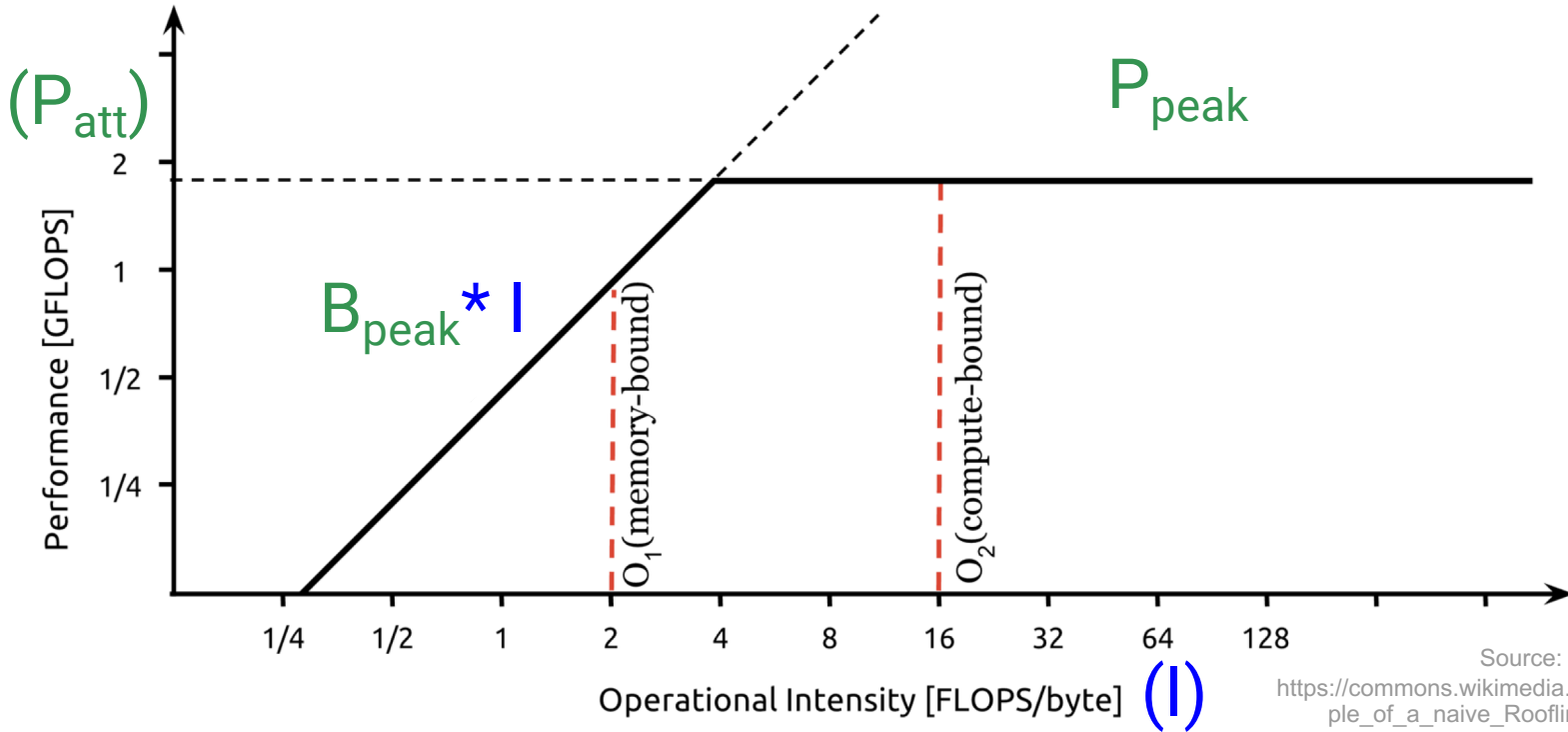
Motivation: Mobile SoCs, Usecases, & IP Blocks

Gables: A Roofline Model for Mobile SoCs

Gables 2-IP Example

Wrap Up

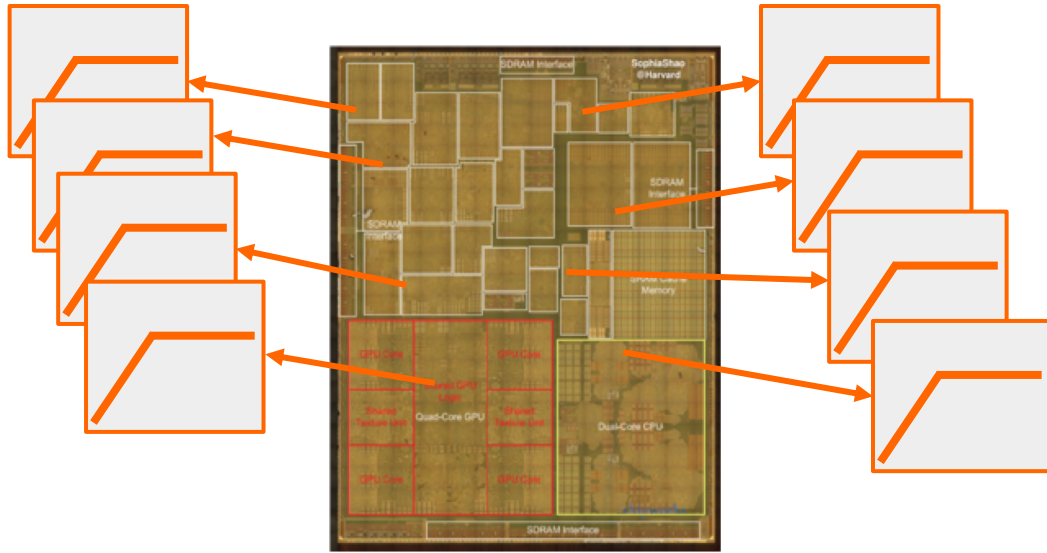
Williams et al., Roofline, CACM 4/2009



Source:
https://commons.wikimedia.org/wiki/File:Example_of_a_naive_Roofline_model.svg

$$P_{att} = \text{MIN}(B_{peak} * I, P_{peak})$$

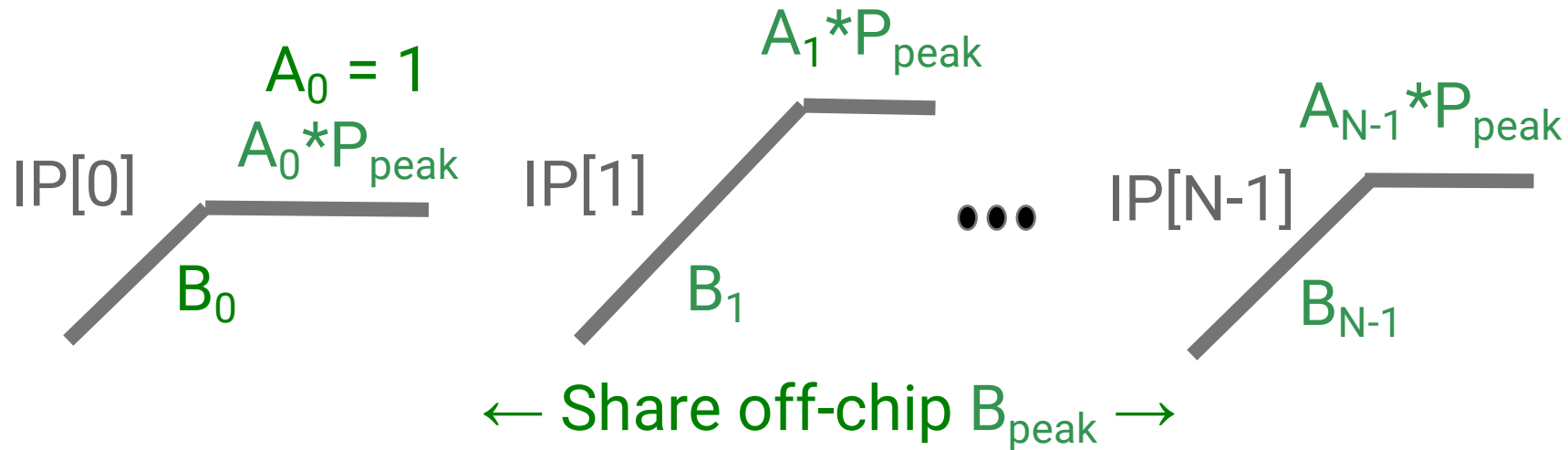
Consumer System on Chip (SoC) & Gables



1. Include Accelerator IP[i]?
2. IP[i] over-provisioned?
3. IP[i] over-communicates?

Gables provides first answer!

Gables for N IP SoC



Usecase at each each IP[i]

- Parallel non-negative work f_i (f_i 's sum to 1)
- Operational intensity l_i operations/byte

Gables Math: Roofline / Work Fraction



$$\text{Roofline: } \begin{array}{l} \text{MIN}(B * l, P_{\text{peak}}) \\ \text{MIN}(B * l, 1 * P_{\text{peak}}) / 1 \end{array}$$

$$1 / T_{\text{IP}[i]} = \text{MIN}(B_i * l_i, A_i * P_{\text{peak}}) / f_i \quad f_i \neq 0$$

$$1 / T_{\text{memory}} = B_{\text{peak}} * l_{\text{avg}} \quad l_{\text{avg}} = 1 / \sum_{i=1, N-1} (f_i / l_i)$$

$$\text{Perf} = \text{MIN}(1/T_{\text{IP}[0]}, \dots, 1/T_{\text{IP}[N-1]}, 1/T_{\text{memory}})$$

Outline

Motivation: Mobile SoCs, Usecases, & IP Blocks

Gables: A Roofline Model for Mobile SoCs

Gables 2-IP Example

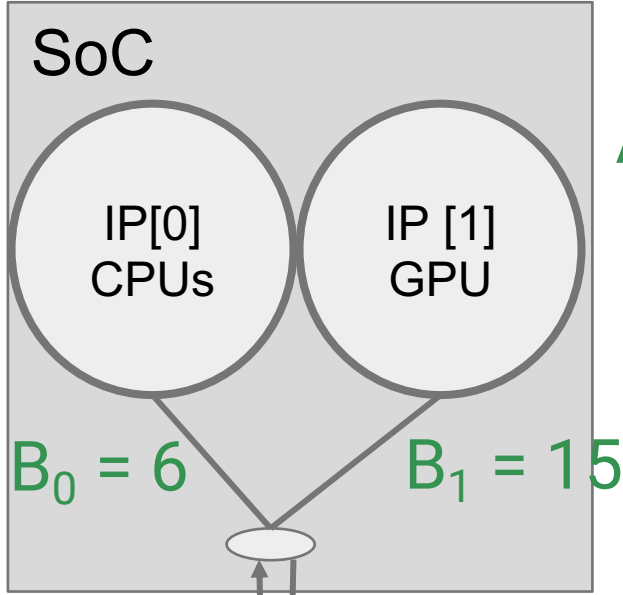
Wrap Up

Gables Example with 2 IP SoC System

$$P_{\text{peak}} = 40$$

$$A * P_{\text{peak}} = 5 * 40 = 200$$

$$B_{\text{peak}} = 10$$



Workload:

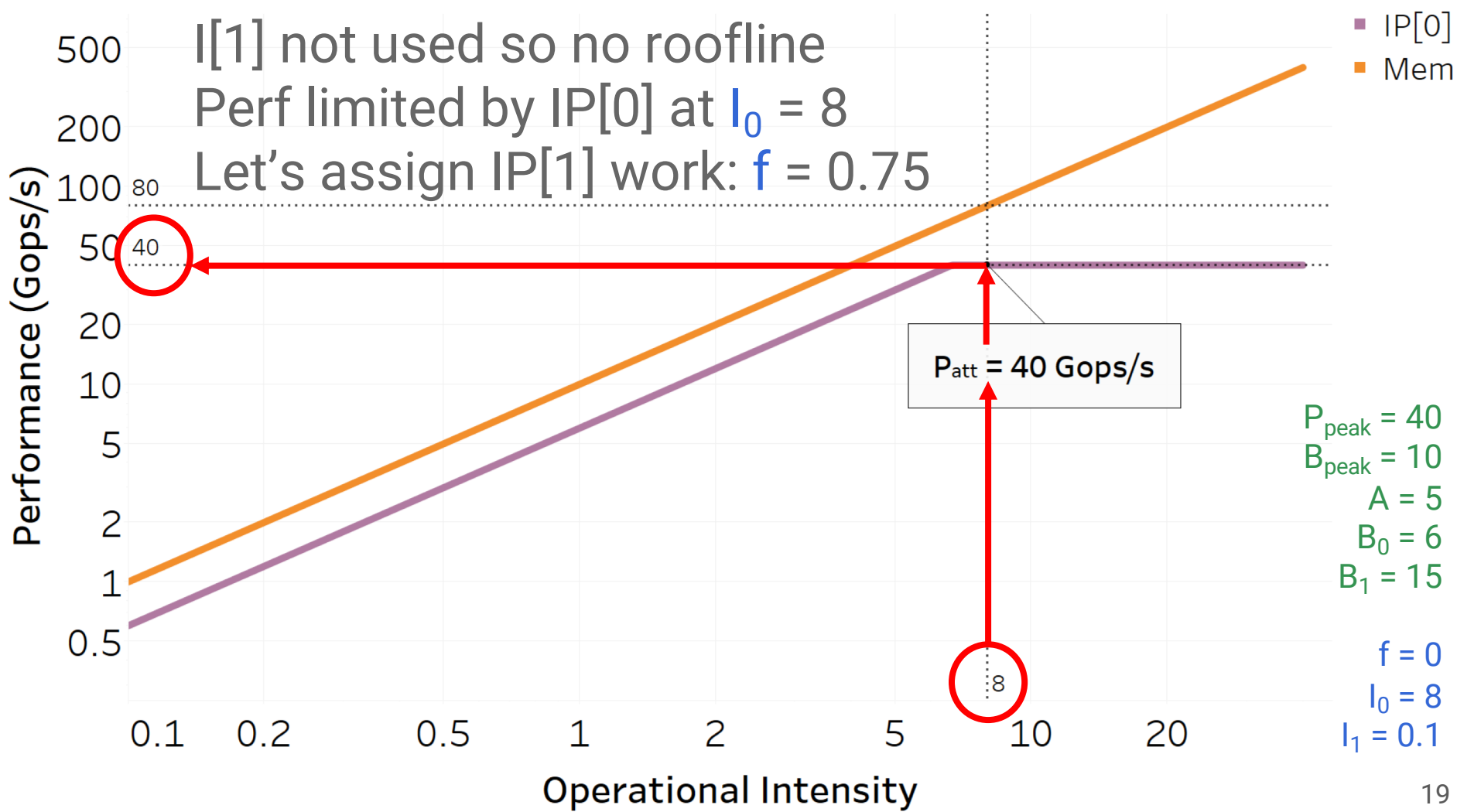
$$f = 0$$

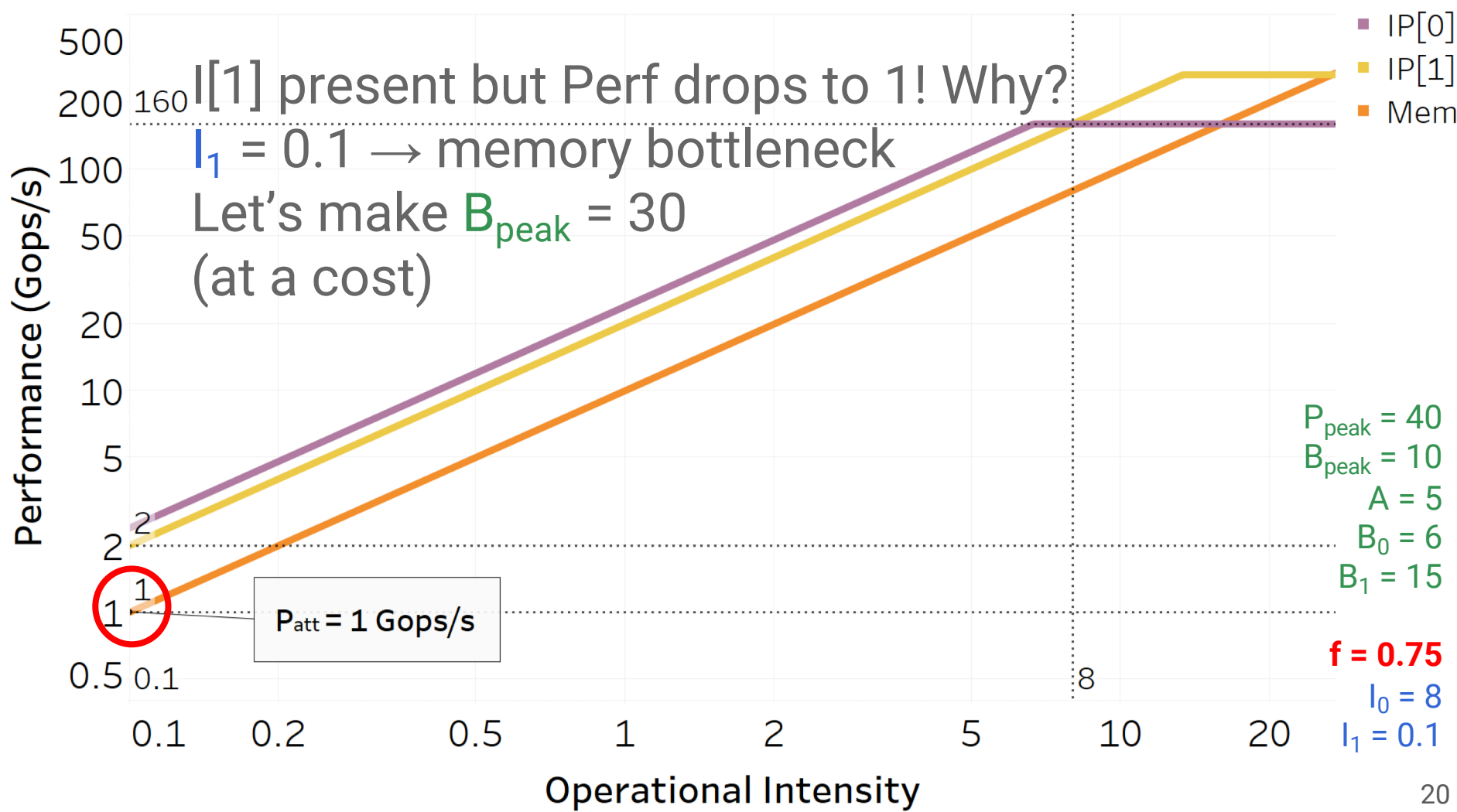
$l_0 = 8 = \text{good caching}$

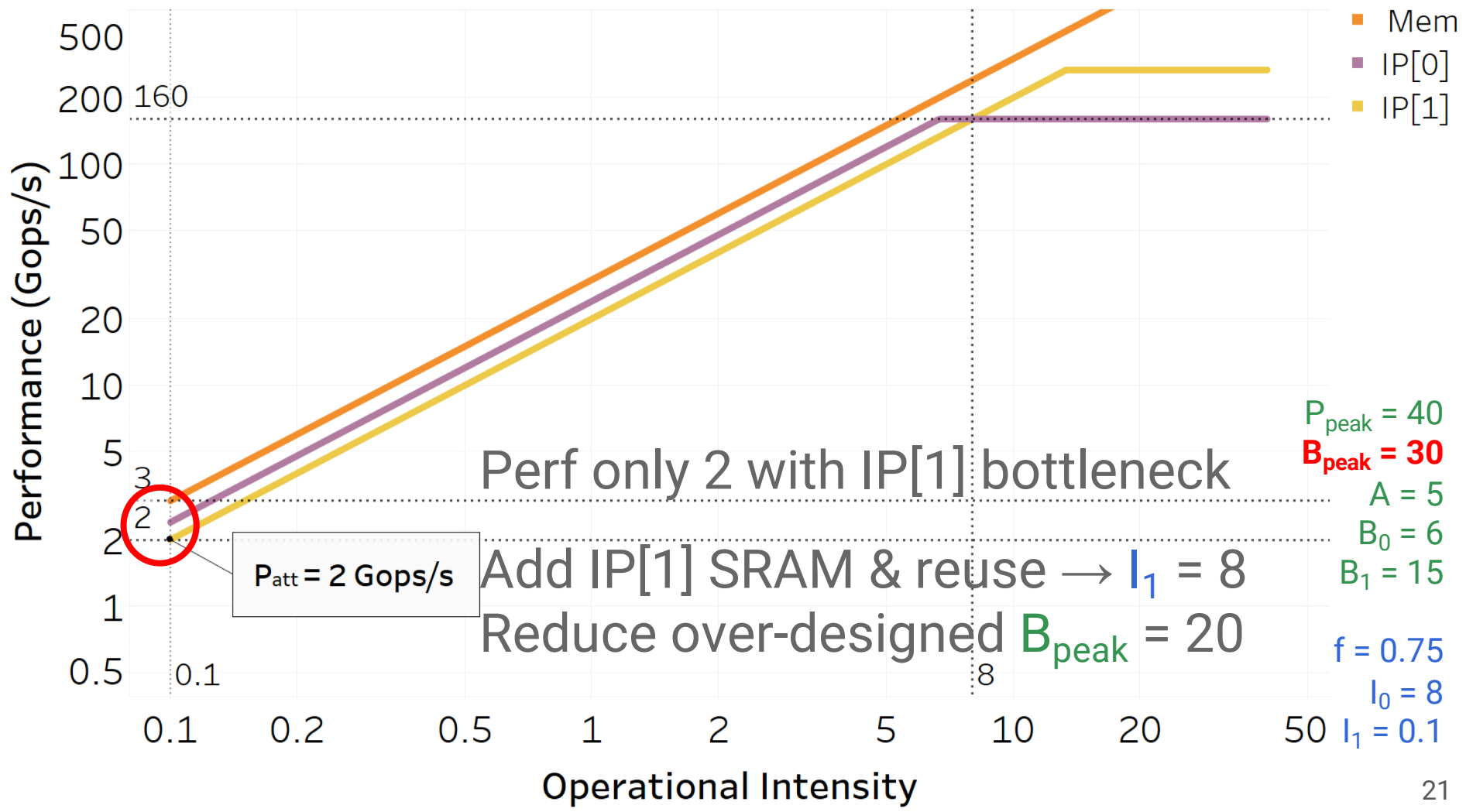
$l_1 = 0.1 = \text{latency tolerant}$

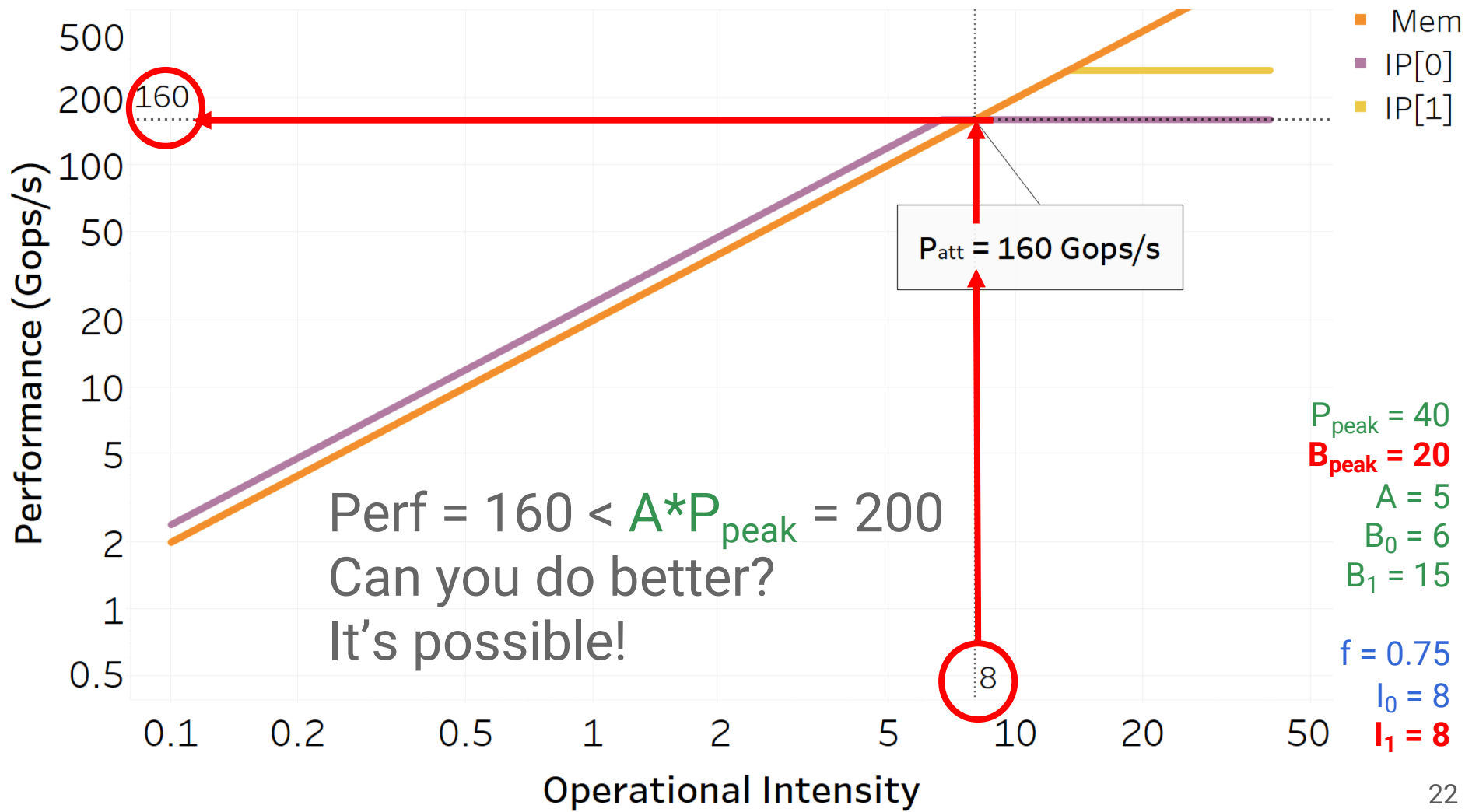
Performance?

Go to research.cs.wisc.edu/multifacet/gables/

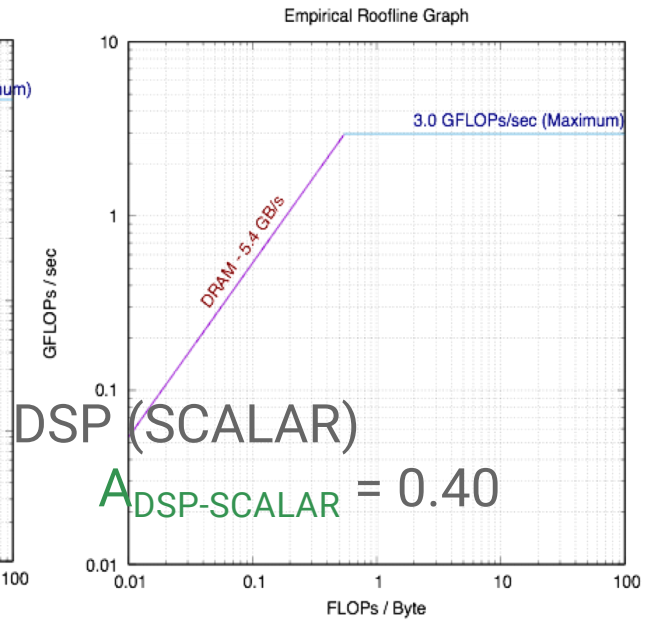
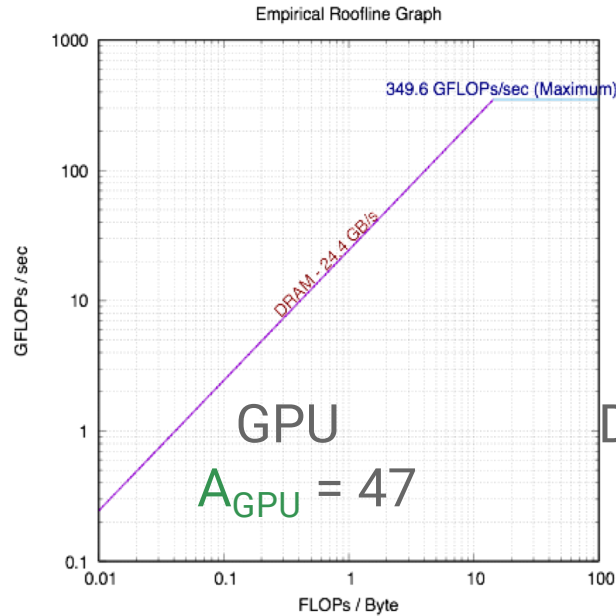
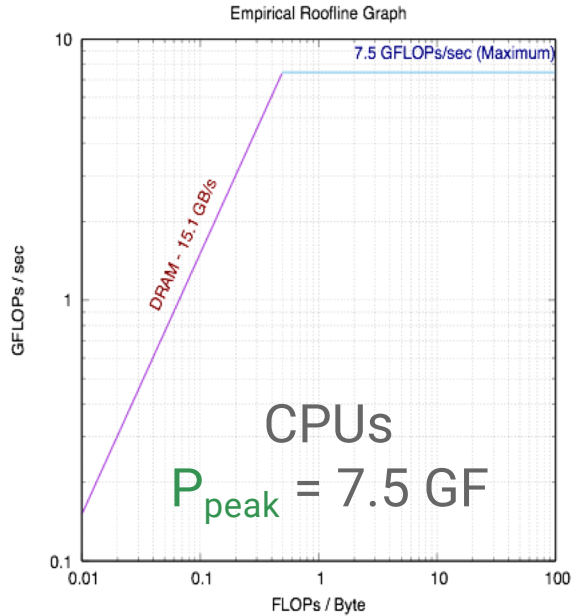








μ Benchmark w/ Qualcomm Snapdragon™ 835



- All elements load from array & vary FP SP op intensity
- Finds empirical lower bound on rooflines

Gables Paper & Home Page

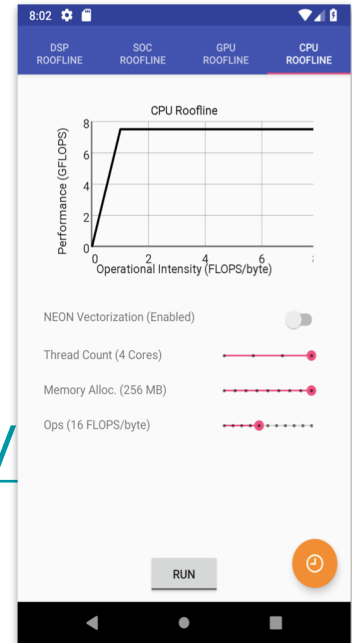
Extensions: memory-side buffer, interconnect, serial work

Interactive tool for 2-IP & 3-IP SoCs

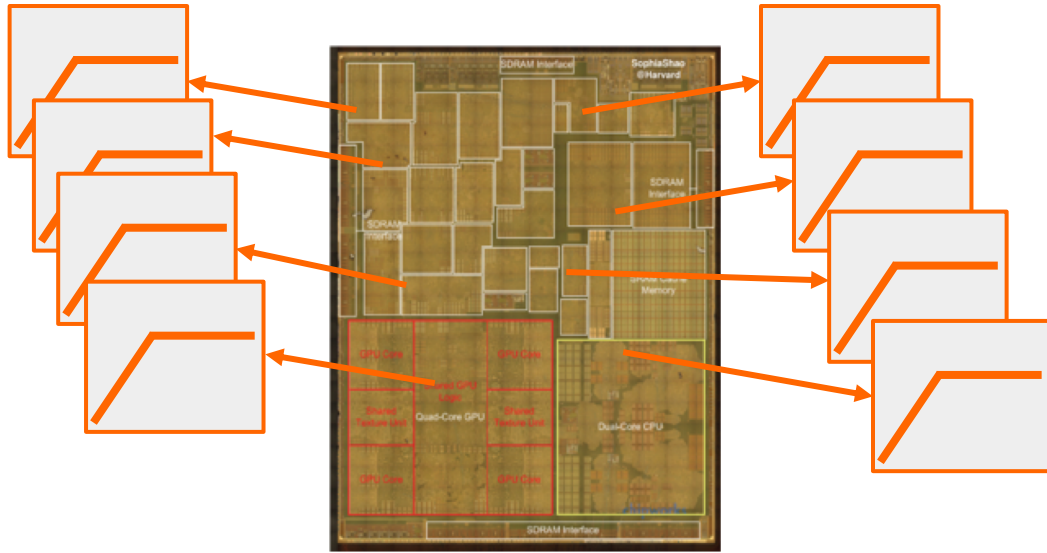
Gables: Open-source Android App



<http://research.cs.wisc.edu/multifacet/gables/>
<https://github.com/harvard-edge/Gables>



Consumer System on Chip (SoC) & Gables



1. Include Accelerator IP[i]? Or give work to enhanced CPUs
2. IP[i] over-provisioned? Make IP[i] acceleration less
3. IP[i] over-communicates? IP[i] less compute; more SRAM

Conjectures

1. Gables is useful for early Mobile SoC planning
1. Valuable to scrutinize each IP's Acceleration & BW
1. Estimating work fraction for “Goldilocks” IP design
1. **Operation intensity** illuminates IP memory reuse