

The ARM logo is displayed in a white, lowercase, sans-serif font. It is positioned on the left side of the slide, set against a background of a blue-tinted, high-angle view of a city at night with glowing lights and a grid of small white plus signs.

arm

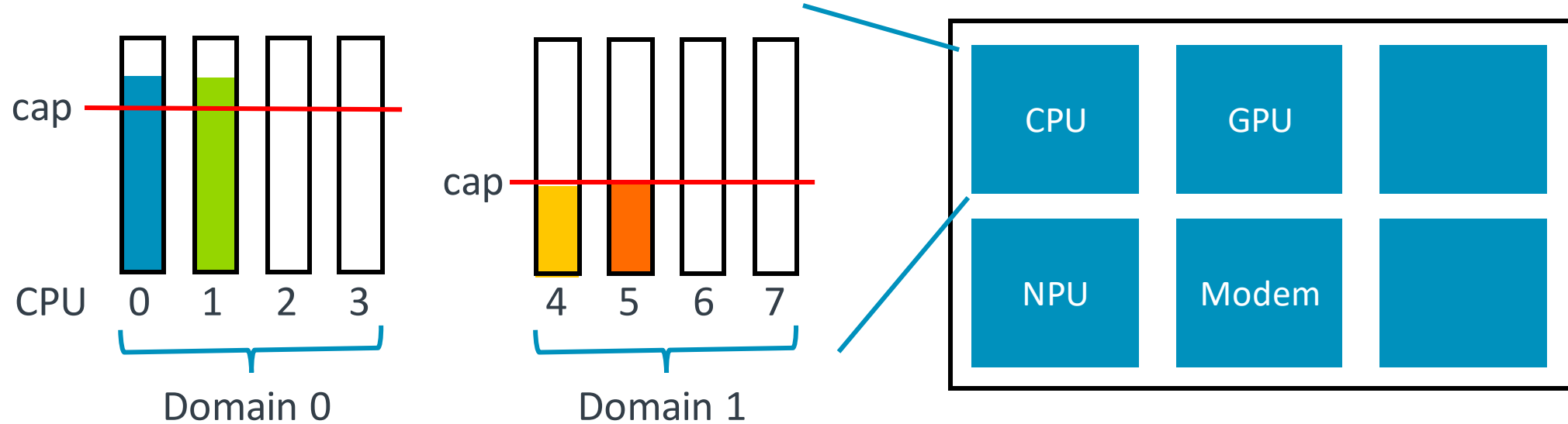
# Performance guarantees under thermal pressure

Morten Rasmussen <[morten.rasmussen@arm.com](mailto:morten.rasmussen@arm.com)>

Linux Plumbers Conference 2019, 9-11 September, Lisbon

# Unpredictable compute bandwidth

- Compute bandwidth is increasingly becoming unpredictable.
- In mobile systems performance capping is a common scenario.
- User-space has no information about minimum compute bandwidth.



# Is best effort compute enough?

- Current Linux kernel model:
  - User creates tasks, the kernel and platform delivers as good performance as it can.
  - User can tweak task placement and cpufreq governors (including util\_clamp) but kernel and/or platform can override most of it.
- SCHED\_DEADLINE is reservation-based and implies a bandwidth guarantee.
- Reservations are not cleared with thermal framework and could be impossible to fulfill.
- Should we have a “guaranteed” performance level that SCHED\_DEADLINE could use for admission control?
- What level of “guarantee” should it provide?
- Who should provide it? DT, ACPI?



The ARM logo is displayed in a white, lowercase, sans-serif font. It is positioned on the left side of the slide, set against a background of a blue-tinted, high-angle view of a city at night with glowing lights and a grid of small white plus signs.

arm

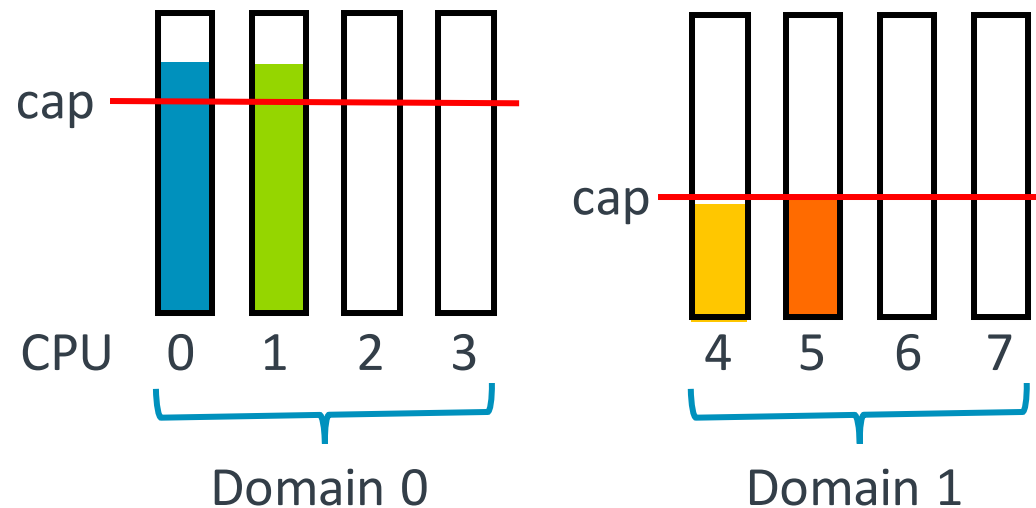
# Task-centric thermal management

Morten Rasmussen <[morten.rasmussen@arm.com](mailto:morten.rasmussen@arm.com)>

Linux Plumbers Conference 2019, 9-11 September, Lisbon

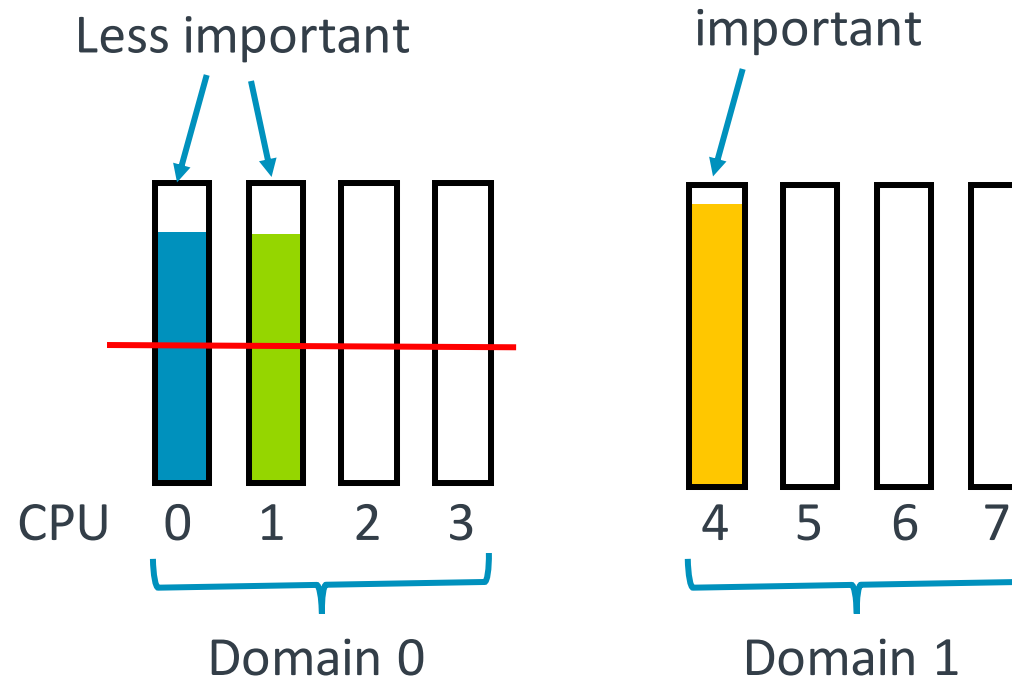
# CPU thermal management in Linux

- Most systems control thermally unsustainable compute demand by performance capping.
- Control is performance domain-centric (clock domain) capping one or more CPUs together.
  - No knowledge about relative importance of tasks on the CPUs in each performance domain.
  - No knowledge about how to best spend the thermal budget.
  - Scheduler might help tasks to “escape” capped CPUs moving the problem somewhere else.



# Task-centric thermal management

- Ideally thermal management should maintain the budget while minimizing the perceived performance impact.
- Middleware/application is best positioned to decide the budget split.
- If tasks would self-adapt their compute demand, it would be even better.



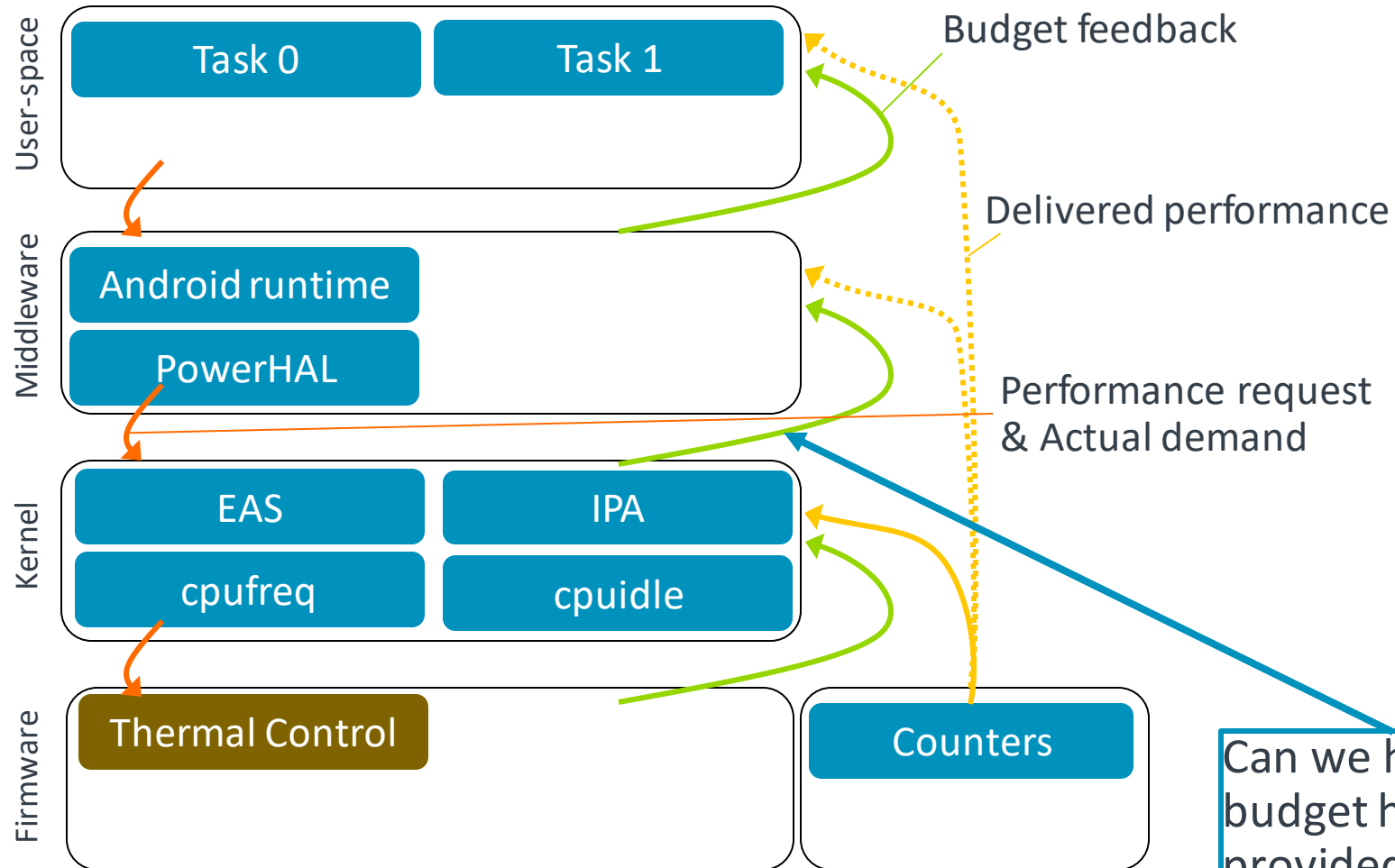
# Power management mechanisms

- Applications
    - Self-adapting apps<sup>1</sup>, util\_clamp capping.
  - Middleware:
    - cgroup bandwidth controller, (util\_clamp capping).
  - Kernel:
    - IPA, cpufreq (DVFS), ~~hot-plugging~~, (cgroup bandwidth controller).
  - Firmware:
    - Frequency capping.
- Controlling compute demand
- Controlling compute bandwidth
- Adapting demand is preferable to reducing compute bandwidth.
  - Middleware or Kernel steps in controlling the bandwidth if applications don't behave.

<sup>1</sup> <https://blogs.unity3d.com/2019/04/01/higher-fidelity-and-smoother-frame-rates-with-adaptive-performance/>



# Performance management hierarchy



Can we have a generic thermal budget headroom metric provided by the thermal framework?