

## Speedup and Power Scaling Models for Heterogeneous Many-Core Systems\*

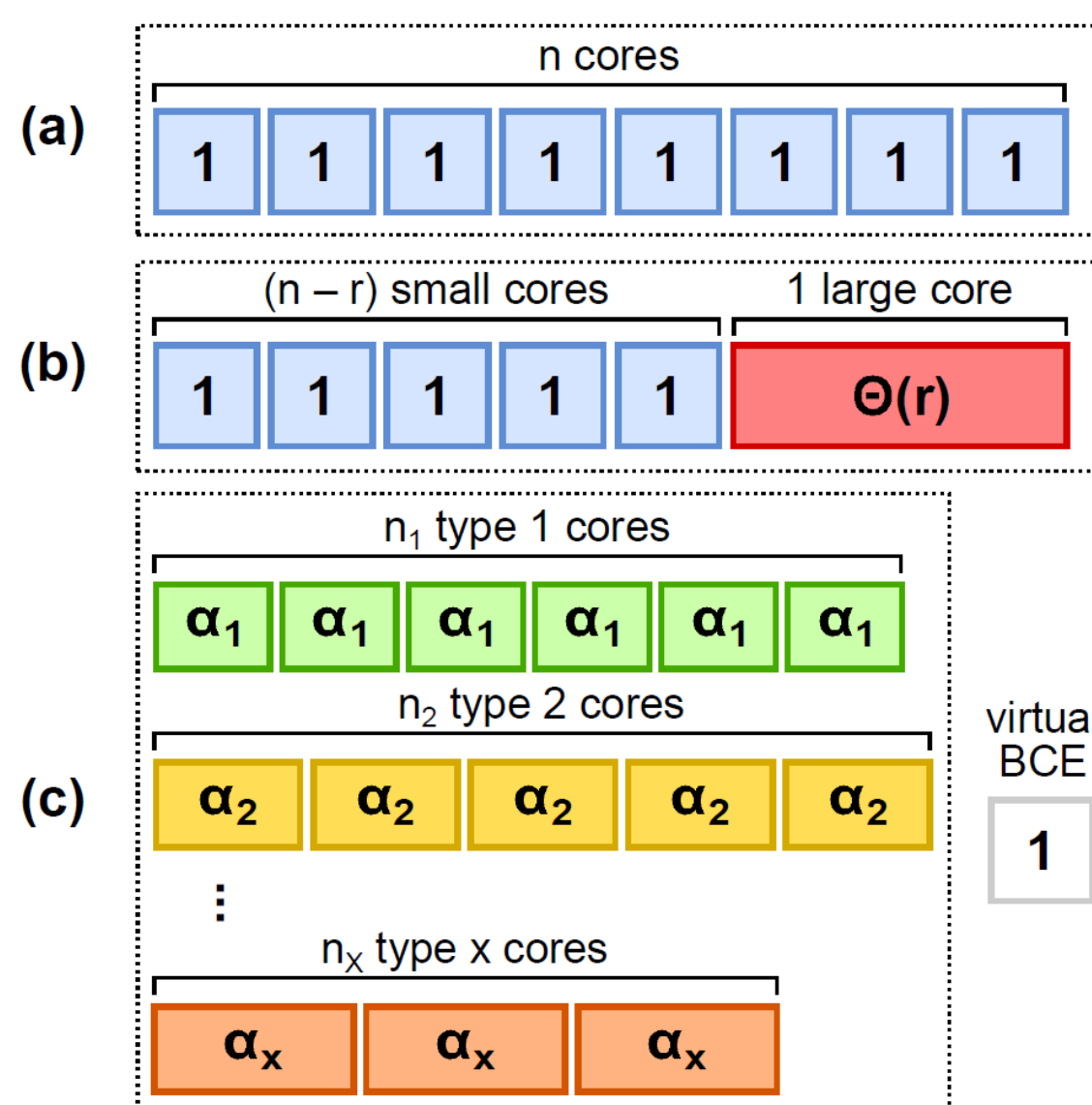
A. Rafiev, M. Al-hayanni, F. Xia, R. Shafik, A. Romanovsky, A. Yakovlev

### INTRODUCTION

- Traditional speedup models help the research community and industry better understand system performance capabilities and application parallelizability.
- We introduce normal form heterogeneity, that supports a wide range of heterogeneous many-core architectures.
- The modelling method aims to predict system power efficiency and performance ranges.
- The models were validated through extensive experimentation on the off-the-shelf big.LITTLE heterogeneous platform and a dual-GPU laptop
- A quantitative efficiency analysis targeting the system load balancer on the Odroid XU3 platform was used to demonstrate the practical use of the method.

\* Accepted for publication in IEEE Transactions on Multi-Scale Computing Systems.

### HETEROGENEITY



- (a) Homogeneous system (classical Amdahl's Law)
- (b) Simple heterogeneous model (Hill-Marty) consisting of 1 big and many little cores.
- (c) **Proposed model:** x types of cores represented by their relative performances.

### AMDAHL'S LAW

**Homogeneous:**

$$S(n) = \frac{1}{(1-p) + \frac{p}{n}}$$

$p$  – parallelization factor,  
 $n$  – number of cores.

**Heterogeneous:**

$$S(\bar{n}) = \frac{1}{\frac{(1-p)}{\alpha_s} + \frac{p}{N_\alpha}}$$

$\alpha_s$  – sequential core performance,  
 $N_\alpha$  – relative performance of all parallel cores.

### WORKLOAD SCALING

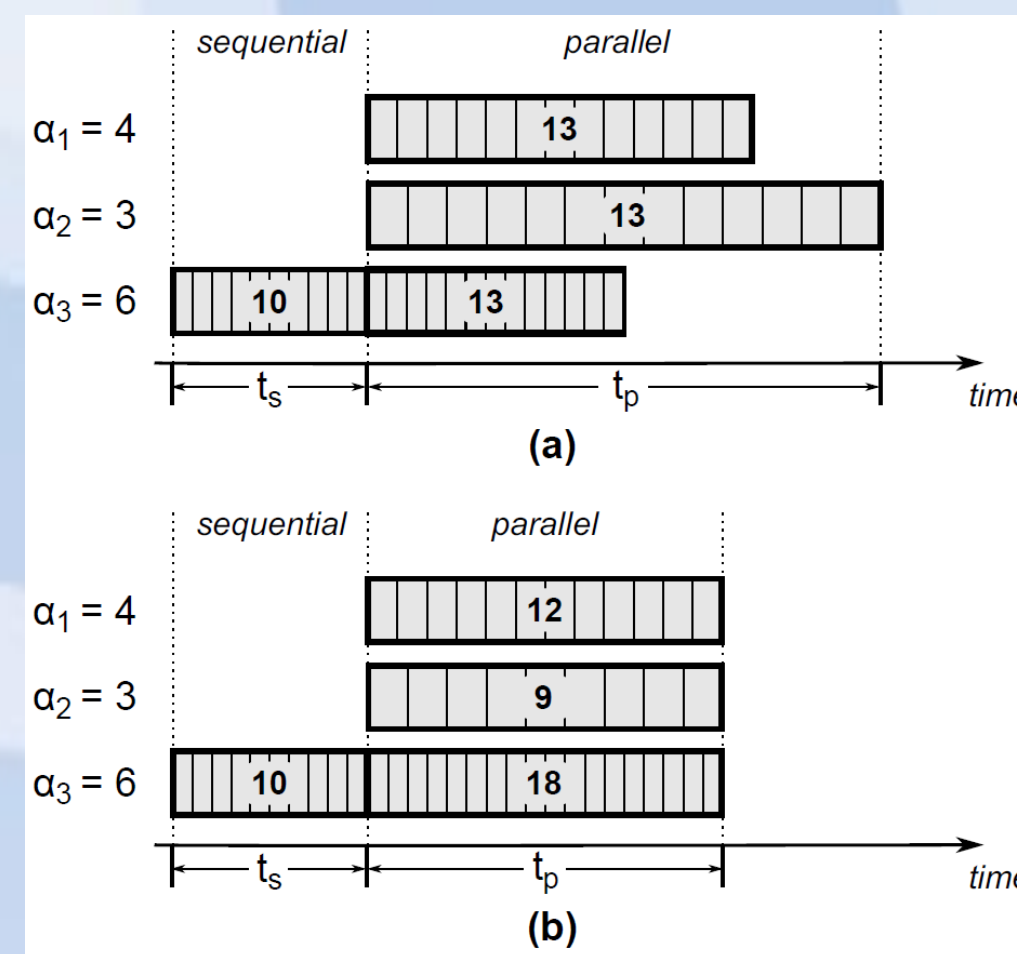
$$I' = h(\bar{n}) \cdot ((1-p)I + pg(\bar{n})I)$$

$I$  – original workload,  $I'$  – scaled workload,  
 $g(n)$  – parallel scaling,  $h(n)$  – proportional scaling.

**General form speedup model:**

$$S(\bar{n}) = \frac{(1-p) + pg(\bar{n})}{\frac{(1-p)}{\alpha_s} + \frac{pg(\bar{n})}{N_\alpha}}$$

### LOAD DISTRIBUTION



(a) **Equal-share (naive):**

$$N_\alpha = \left( \sum_{i=1}^x n_i \right) \cdot \left( \min_{1 \leq i \leq x} \alpha_i \right).$$

(b) **Balanced (ideal):**

$$N_\alpha = \sum_{i=1}^x \alpha_i n_i.$$

### POWER MODELLING

$$W_{total} = W_0 + W(\bar{n}),$$

where  $W_0$  is **background** power and  $W$  is **effective** power.

For  $w$  – BCE power, and  $(\beta_1, \dots, \beta_x)$  – relative core powers:

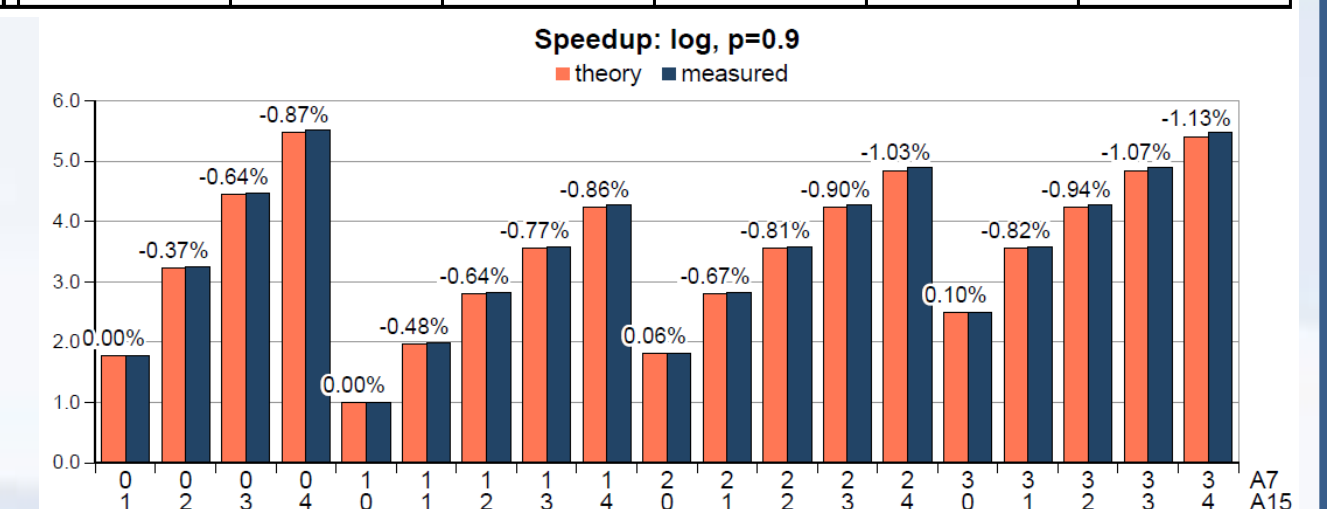
$$W(\bar{n}) = w D_w(\bar{n}) S(\bar{n})$$

$$D_w(\bar{n}) = \frac{\frac{\beta_s}{\alpha_s} (1-p) + pg(\bar{n}) \frac{N_\beta}{N_\alpha}}{(1-p) + pg(\bar{n})}$$

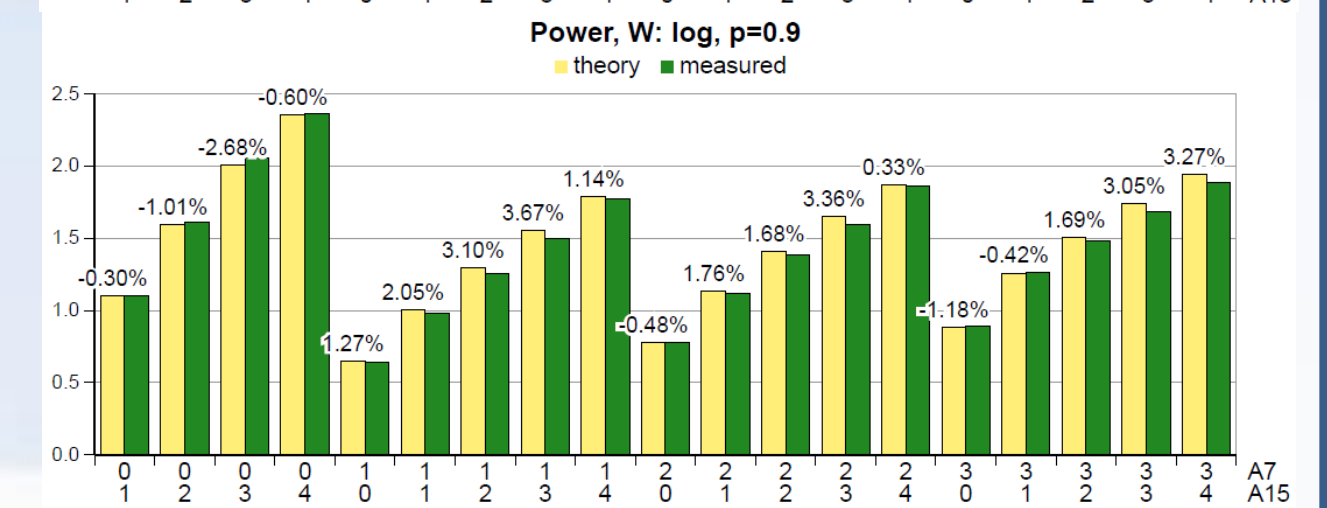
### ODROID XU3

benchmark base workload	sqrt 40000		int 40000		log 40000	
	A7	A15	A7	A15	A7	A15
core type $i$						
measured execution time, ms	49969	53206	52844	42665	41820	23506
measured active power, W	0.2655	0.8361	0.2760	0.8305	0.3036	0.9496
power measurement std dev	0.82%	0.18%	0.96%	0.87%	0.93%	0.42%
calculated effective power, W	0.1158	0.4887	0.1264	0.4830	0.1540	0.6022
$\alpha_i$	1	0.9392	1	1.2386	1	1.7791
$\beta_i$	1	4.2183	1	3.8221	1	3.9094

Speedup error < 1.2%



Power error < 5.6%



### PARSEC BENCHMARKS

Evaluating system load balancer quality:  $N_{low} < N_{meas} < N_{high}$

