

Sharing Incentives and Fair Division for Multiprocessors

Seyed Majid Zahedi and Benjamin C. Lee (Duke University)

1. Case for Sharing

Big Servers

Hardware is under-utilized
Sharing amortized power

Heterogeneous Users

Tasks are diverseUsers are complementary



2. Conventional Wisdom in Computer Architecture

Users must share

Overlooks strategic behavior

Fairness policy is equal slowdown

Users prefer flexibility

Sharing ChallengesAllocate multiple resourcesEnsure fairness

Intel Sandy Bridge E die [www.anandtech.com]

• Fails to encourage envious users to share

Heuristic mechanisms enforce equal slowdown

• Fail to give provable guarantees

3. Resource Elasticity Fairness (REF)

REF is an allocation mechanism that guarantees gametheoretic desiderata for shared chip multiprocessors

Sharing Incentives

Users perform no worse than under equal division

Envy-Free

No user envies another's allocation

Pareto-Efficient

No other allocation improves utility without harming others
 Strategy-Proof
 No user benefits from lying

4. Cobb-Douglas Utility

 $\mathbf{u}(\mathbf{x}) = \prod_{\mathsf{r}=1}^{\mathsf{R}} \mathbf{x}_{\mathsf{r}}^{lpha_{\mathsf{r}}}$

u utility (e.g., performance)

- k_r allocation for resource r (e.g., cache size)
- $\alpha_{\mathbf{r}}$ elasticity for resource \mathbf{r}
- Cobb-Douglas fits preferences in computer architecture
- Exponents model diminishing marginal returns
- Products model substitution effects



6. Mechanism for Resource Elasticity Fairness



7. Dynamic Resource Allocation

8. Allocation of Heterogeneous Cores with

Limitation of prior mechanisms

- Fairness is guaranteed within a single epoch
- Fairness is achieved by sacrificing performance

Dynamic system behavior

- System is defined by sequence of epochs
- Users may enter and leave the system
- Users demand for resources may change
- Resource availability may change

Token System

Virtual currency is used to allocate big cores
Users send a request for big core
Users spend one token for big cores
Users receive redistributed tokens

Users are given the same number of tokensMechanism guarantees big cores for certain number of epochs

Threshold strategy is used to optimize performance
Users request big cores, if small core cannot provide comparable performance