

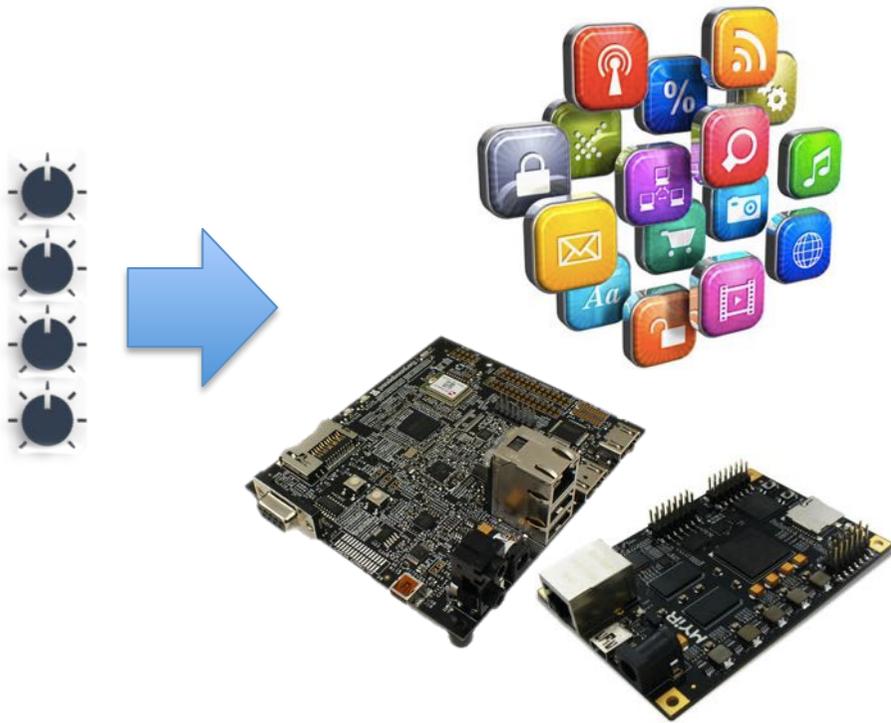


POLITECNICO
MILANO 1863

Energy Efficient Computing Systems Exploiting Online Tuning and Output Quality Management

Gianluca Palermo

Extra-Functional Properties



Functional Description



What to do...

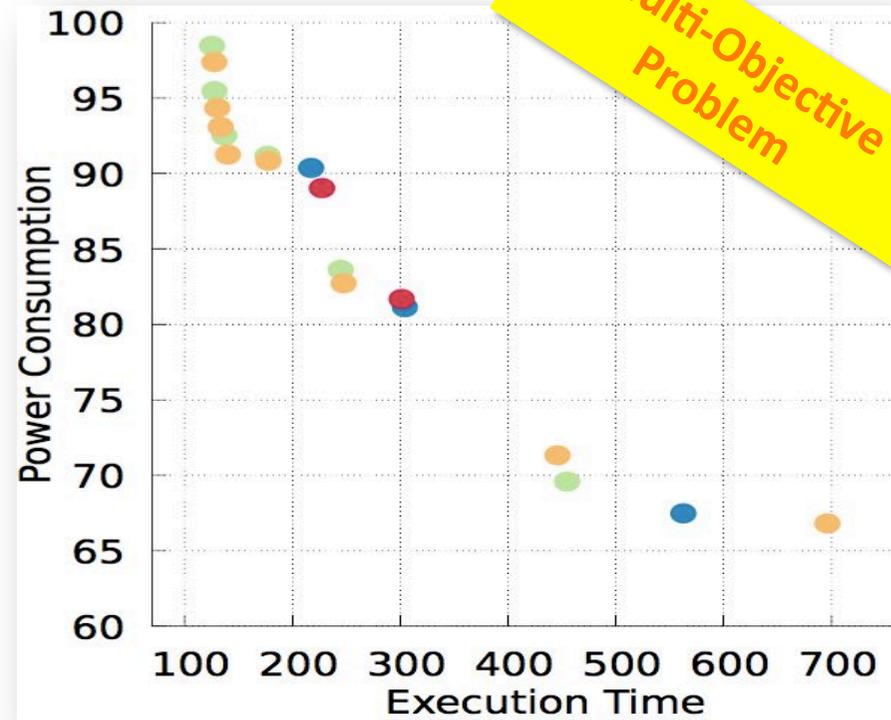
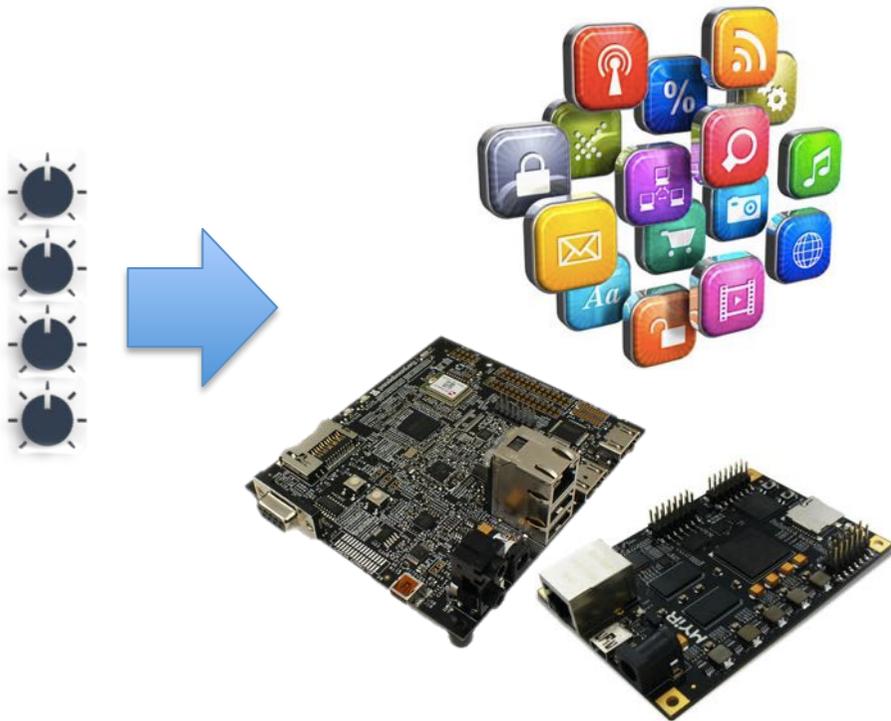


Extra-Functional Properties

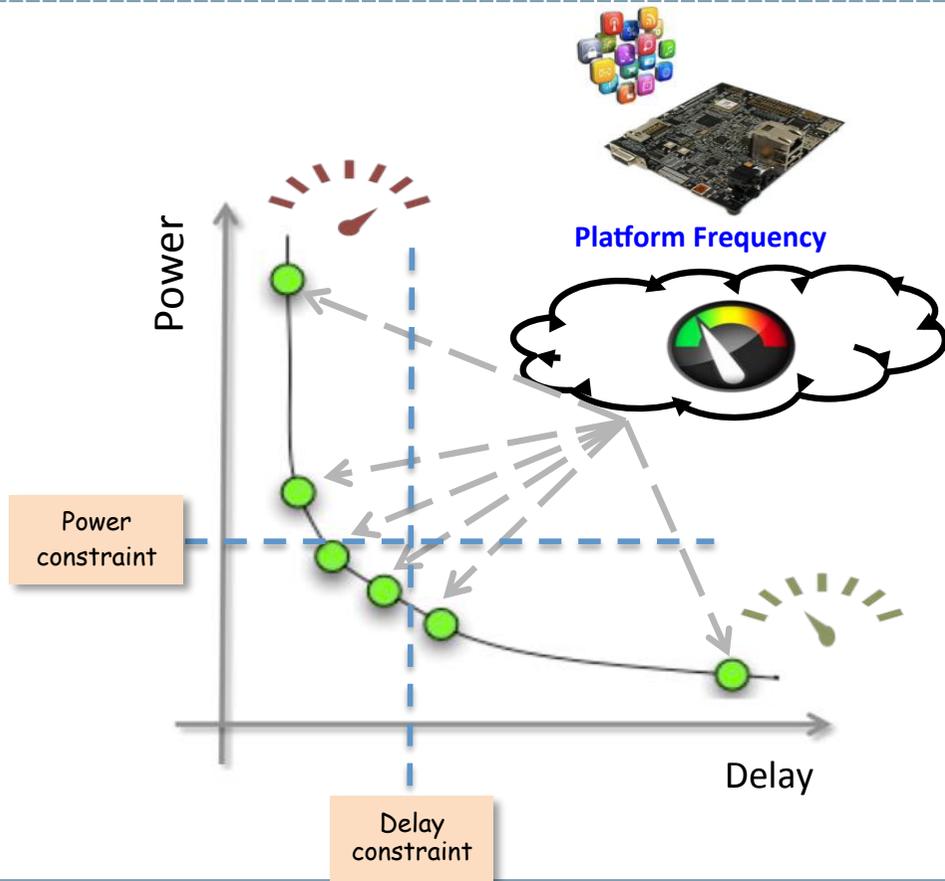


How It is done...

Extra-Functional Properties

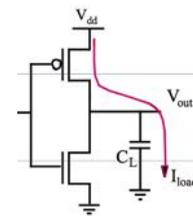


Simple Power-Performance trade-off



Dynamic Power

$$\text{Power} = f * C * V^2 + P_{\text{static}}$$



Energy constraint
???



$$\text{Energy} = \text{Power} * \text{Delay}$$

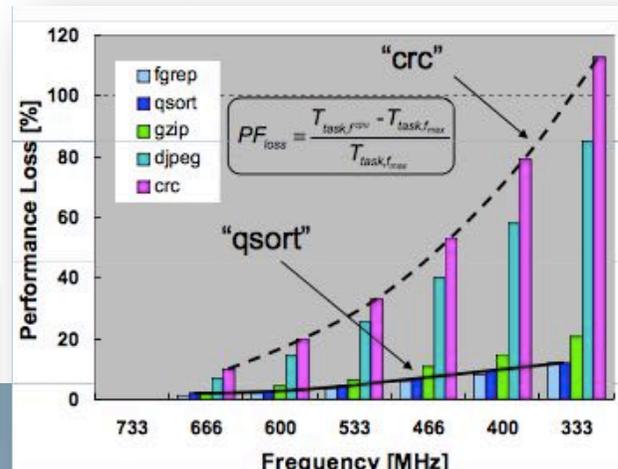
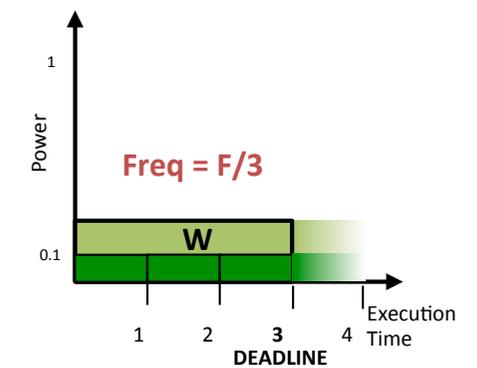
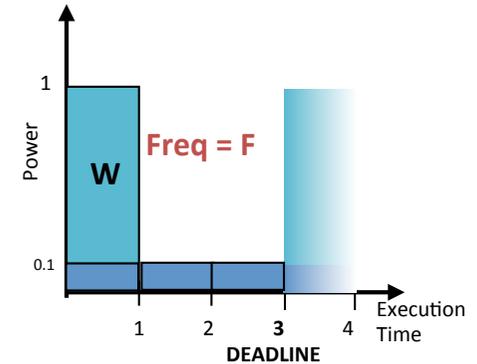
Energy Vs Power

- *Power* poses constraints
 - E.g. power delivery or cooling solution
- *Energy* is in most of the cases the ultimate metric
 - It measures the *cost* of performing a fixed task
- ***How can I reduce Energy?***
 - Voltage and Frequency Scaling
 - Dynamic Power Switching
 - [...]



(Dynamic) Voltage and Frequency Scaling (DVFS)

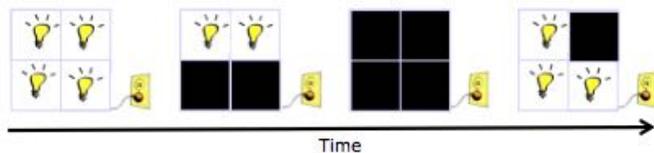
- Small reductions in voltage can be very significant
 - Power is proportional to frequency and square of voltage
- Effective way to reduce energy by providing *just-enough* computing power
- F and V are not independent variables $f \propto (V_{dd} - V_{th})^\alpha / V_{dd}$
- Choosing the right VF operating point is not straightforward



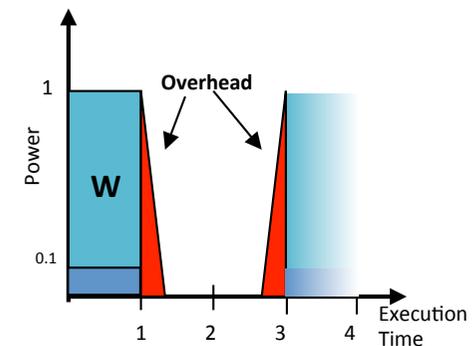
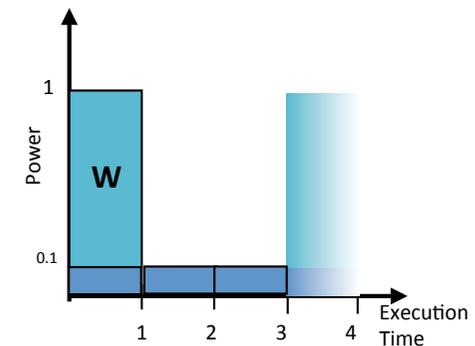
K. Choi, R. Soma, M. Pedram
 "Dynamic voltage and frequency scaling based on workload decomposition" ISLPED 2004

Dynamic Power Switching (DPS)

- Keeping devices powered-on consumes power
 - Cut or reduce power on idle device portions

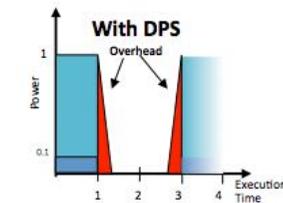
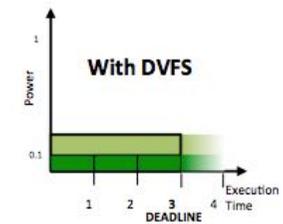
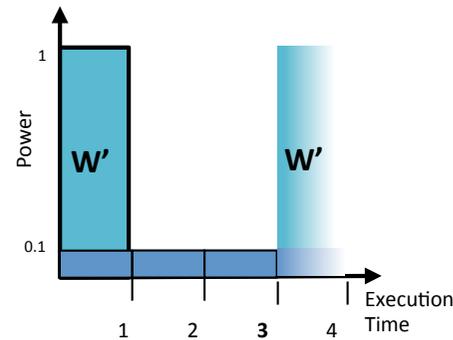
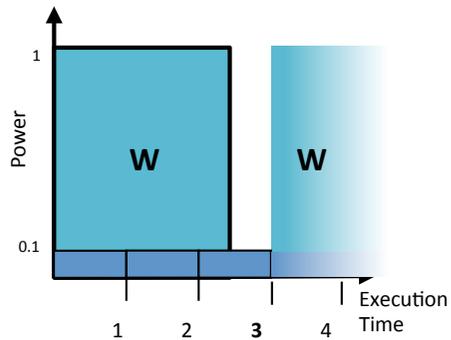


- Significant overheads for switching to deep low-power states
 - Different low-power states
 - Some of them need the context save/restore
 - E.g. Retention states → off states
 - Wakeup latency reduces the responsiveness
- In some contexts also called *Race-to-Halt*



... Longer *Week-Ends* are my Dream

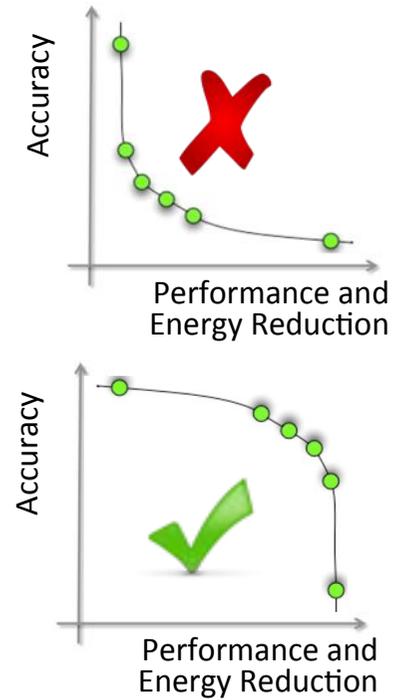
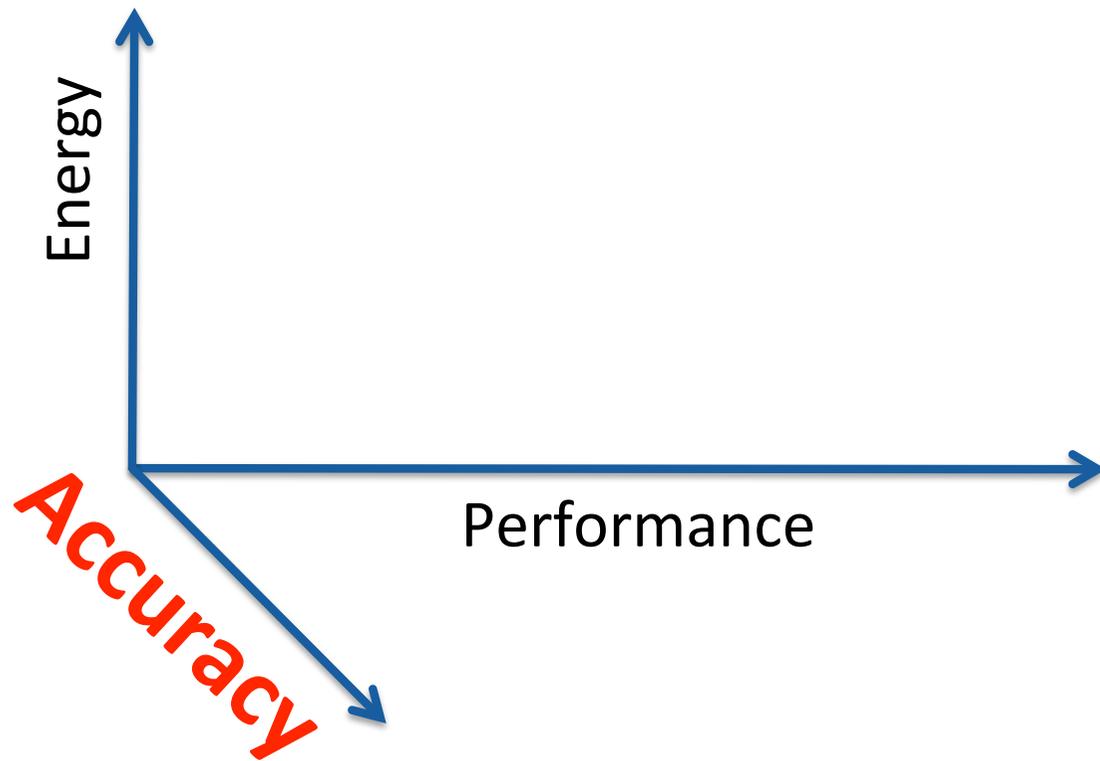
Is there any way to reduce the amount of computation to be performed to consume less?



Approximate Computing



Approximate Computing



Approximate Computing Applications

- Approximation adds **complexity**
 - Not all codes can handle it
- ... but many expensive applications or kernels are naturally error-tolerant
 - Make use of *analog inputs*
 - E.g. operating noisy real-world data from noisy sensor
 - Provide *analog output*
 - E.g. targeting human perception
 - Provide multiple *good-enough* results or *no unique answer*
 - E.g. web search
 - Compute *iteratively towards convergence*
 - E.g. convergent applications over the number of iterations.



V. Chippa et al. "Analysis and characterization of inherent application resilience for approximate computing" DAC 2013.

Possible Application Domains...



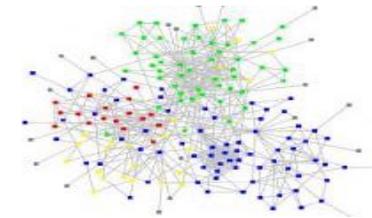
Image Processing



Robotics



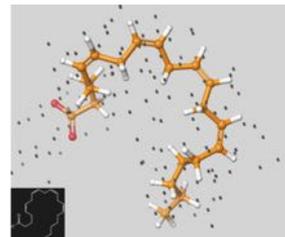
Big Data Analytics



Graph Analytics



Multimedia



Drug Discovery



Traffic Prediction

... where 100% of accuracy not always required

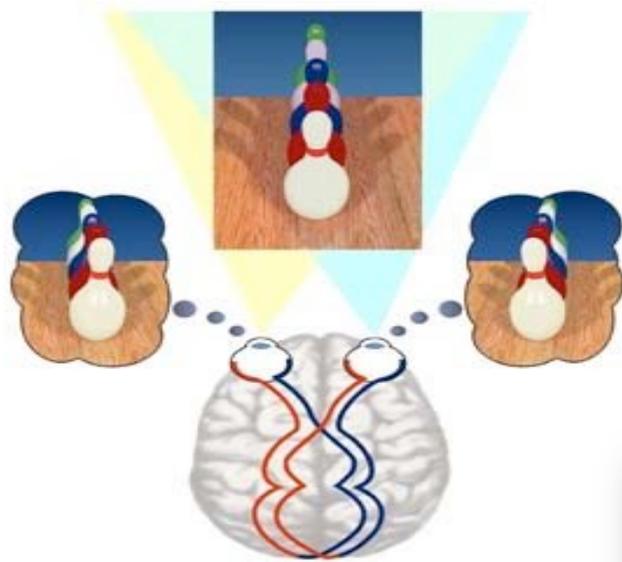
HOWEVER ...

... we have to pay attention



No Traffic

Let's have an example...



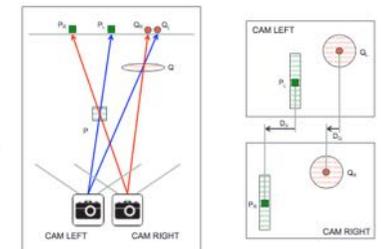
2 eyes = 3 dimensions



Left camera

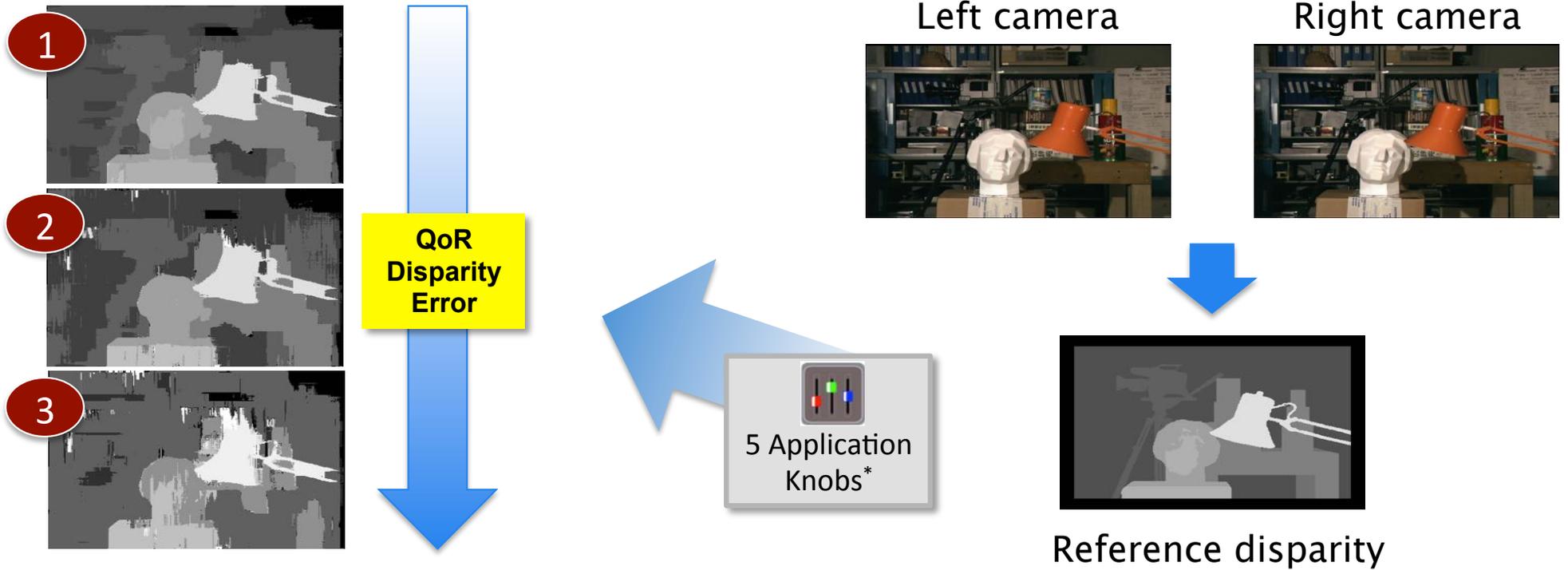


Right camera



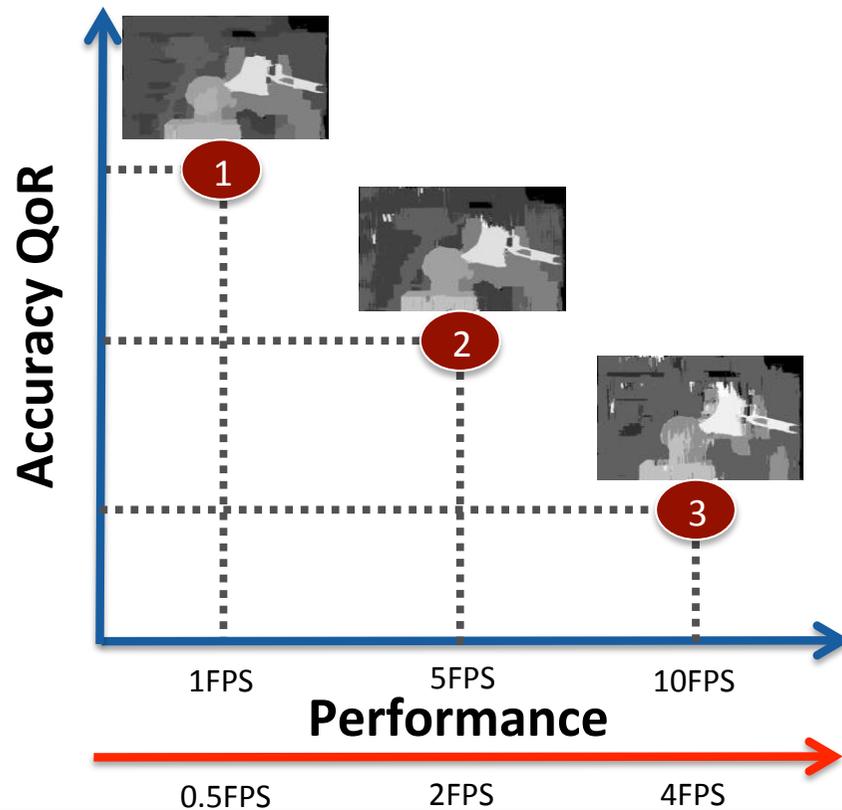
Reference disparity

Tunable Stereo Matching



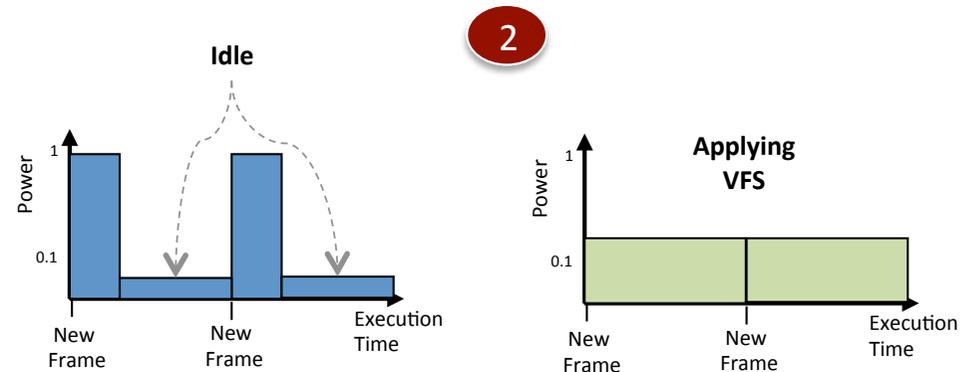
Trading-off Accuracy

D.Gadioli et al. "Application Autotuning to support runtime adaptivity in multicore architectures" SAMOS 2015



Extra-functional requirements:
What if ...

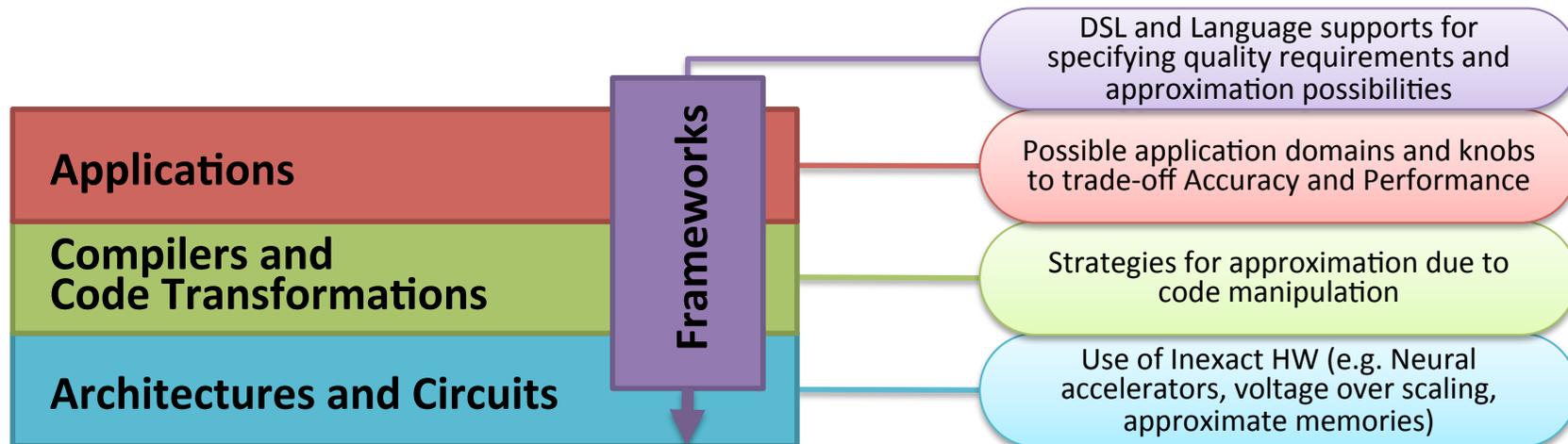
1. Performance = 4FPS 2
2. Performance = F(Speed) 1 2 3
3. Min Energy; QOR > 50%; Perf >= 1



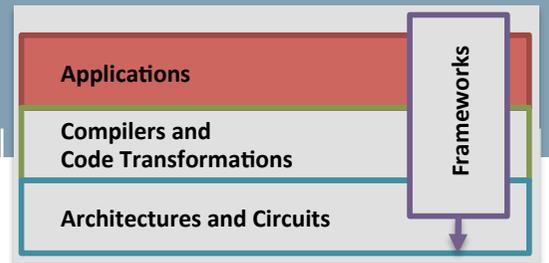
State of the Art...

In literature AC has been faced from different perspectives and at different levels

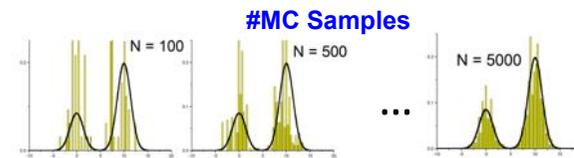
S. Mittal, "A survey of techniques for approximate computing,"
ACM Computing Surveys, pp. 1–34, 2016.



Application-level Approximations



- *Application parameters*



- *Task skipping*

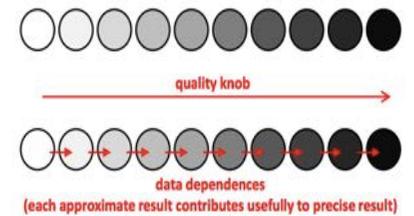
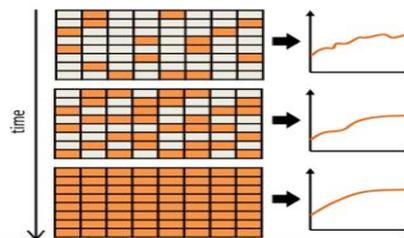
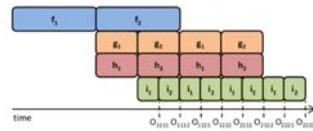
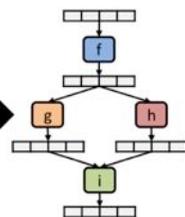


Computation model that represents an approximate application as a pipeline

J. SanMiguel et al "The anytime automaton." ISCA 2016.

```

prologue ();
f ();
g ();
h ();
i ();
epilogue ();
    
```

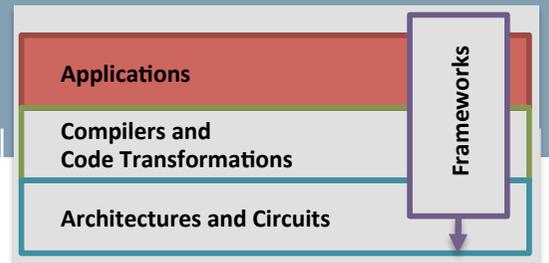


V. Vassiliadis et al " Exploiting Significance of Computations for Energy-Constrained Approximate Computing" IJPP 2016.

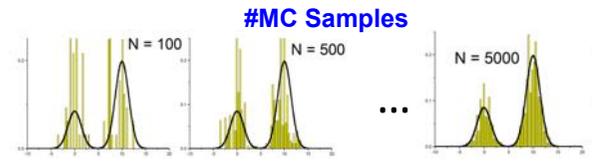
```

#pragma omp task [significant(...)] [label(...)]
[in(...)] [out(...)] [approxfun(function())]
    
```

Application-level Approximations



- *Application parameters*



- *Task skipping*



- *Multiversioning*

- Considering different performance/accuracy trade-offs

DSL or annotation based approaches

J Ansel et al. "PetaBricks: a language and compiler for algorithmic choice" PLDI 2009

B. Woongki et al "Green: a framework for supporting energy-conscious programming using controlled approximation." PLDI 2010

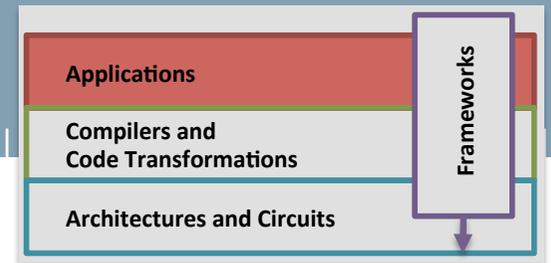
```

1 transform Sort
2 from A[n]
3 to B[n]
4 {
5   from(A a) to(B b) {
6     tunable WAYS;
7     /* Mergesort */
8   } or {
9     /* Insertionsort */
10  } or {
11   /* Radixsort */
12  } or {
13   /* Quicksort */
14  }
15 }
    
```

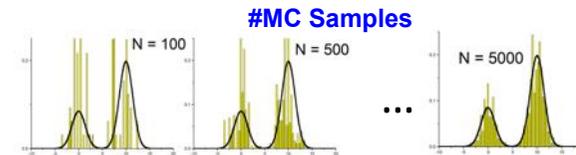
```

switch (get_quality_level(/*...*/)) {
case 0: foo(...) break;
case 1: foo_approx_Q1(...) break;
case 2: foo_approx_Q2(...) break;
/* ... */
}
    
```

Application-level Approximations



- *Application parameters*



- *Task skipping*



- *Multiversioning*

- Considering different performance/accuracy trade-offs

- *Approximate Memoization*

- User partial key for the lookup

(p && 0xffff0000)

```

switch (get_quality_level(/*...*/) {
  case 0: foo(...) break;
  case 1: foo_approx_Q1(...) break;
}

float foo (float p) {
  /* code without side effects */
}

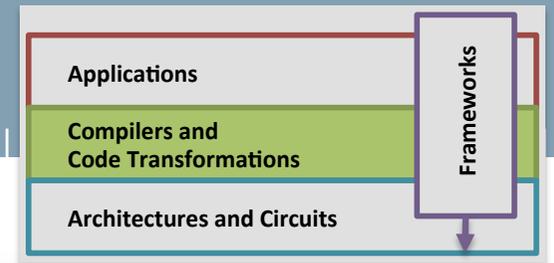
float foo_wrapper(float p){
  float r;
  /* already in the table ? */
  if (lookup_table(p, &r)) return r;
  /* calling the original function */
  else r = foo(p);
  /* updating the table or not */
  update_table(p, r);
  return r;
}
    
```

M. Alvarez et al "Fuzzy Memoization for Floating-Point Multimedia Applications" TCOM 2005.

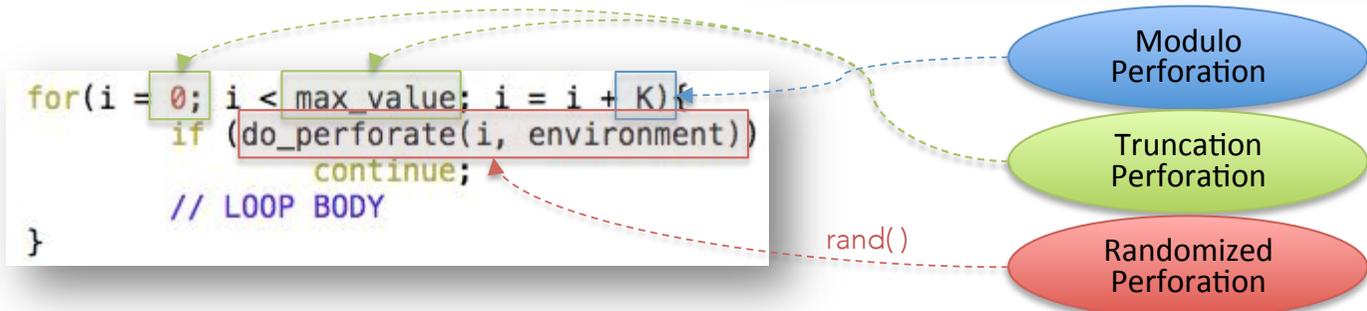
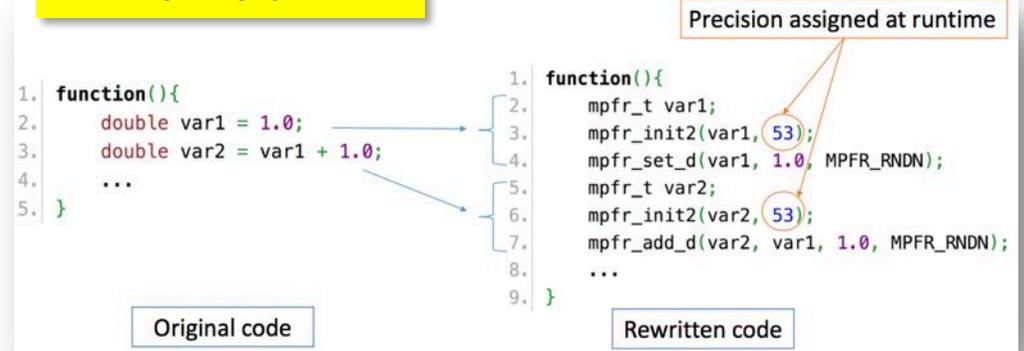
Compiler Approaches

- Precision Scaling
 - FP64->FP32->FP16
 - Float2Int
 - Custom precision

- Loop Perforation

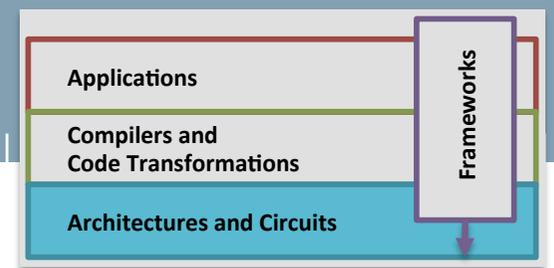


N. Ho et al. "Efficient floating point precision tuning for approximate computing," ASP-DAC 2017.

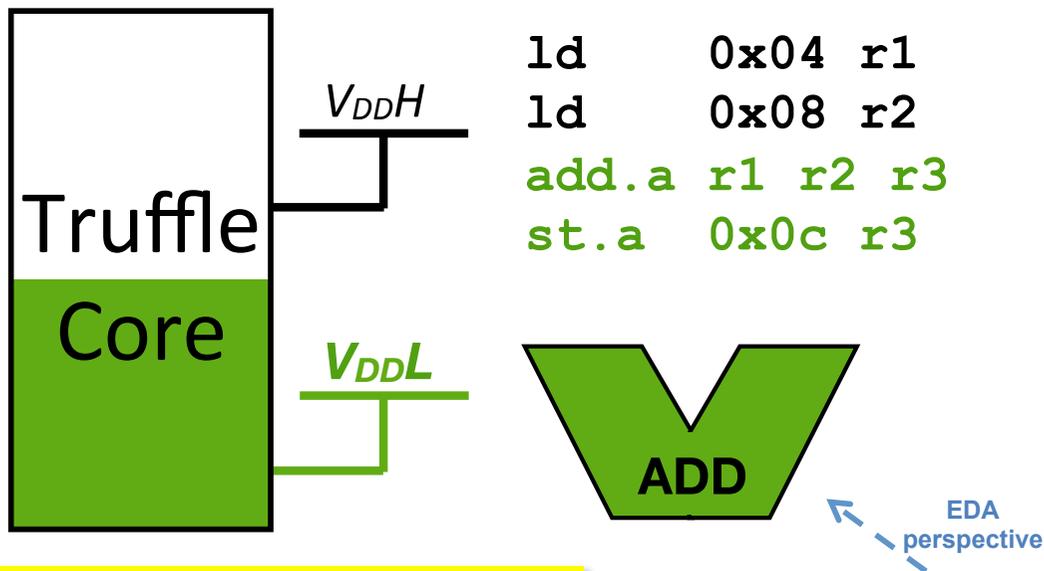


S. Misailovic et al. "Quality of service profiling." ICSE 2010

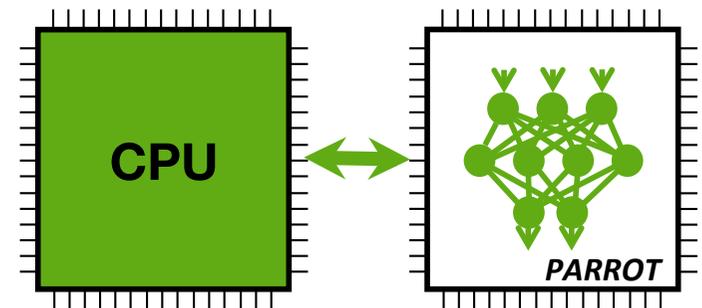
Approximate Hardware



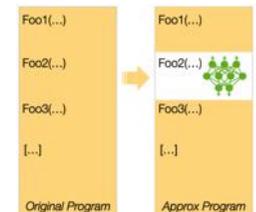
Very large literature ... Dual VDD architectures, Approximate Adders/Multipliers, NPU etc...



H. Esmaelizadeh et al. "Neural Acceleration for General-Purpose Approximate Programs" MICRO 2012



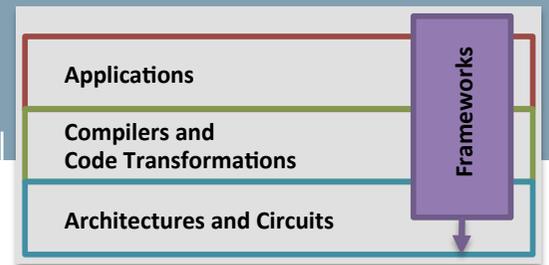
- Efficient hardware implementations
- Capable to mimic many **APPROXIMABLE** computations (after training)
- Fault tolerant



H. Esmaelizadeh et al. "Architecture support for disciplined approximate programming" ASPLOS 2012

S. Venkataramani et al, "SALSA: systematic logic synthesis of approximate circuits" DAC 2012
I. Scarabottolo et al "Circuit Carving: A Methodology for the Design of Approximate Hardware" DATE 2018

Approximate Frameworks (1)



EnerJ

- Java extension with for Approximate Computation

```
Construct
@Approx, @Precise, @Top
endorse(e)
@Approximable
@Context
_APPROX
```

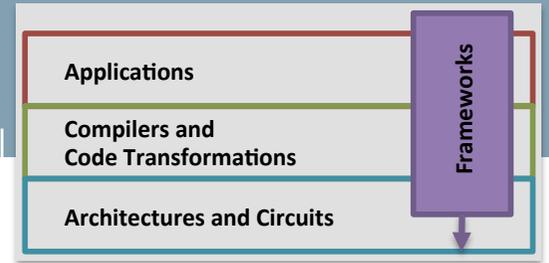
```
@approx    int a = ...;
@precise   int p = ...;
```

- Including Compiler and Runtime envisioning the usage of approximate HW

```
class FloatSet {
    @context float[] nums = ...;
    float mean() {
        calculate mean
    }
    @approx float mean_APPROX() {
        calculate mean of half
    }
}

@approx FloatSet someSet = ...;
someSet.mean(); ←
```

Approximate Frameworks (2)



Functional Description



C/C++ w/ OpenMP, MPI, OpenCL, Matlab



Extra-Functional Requirements and Approximation aware-transformations
E.g. Adaptivity features, Tuning Knobs, Code transformations, Power/Performance constraints



Multiversioning aspect including Modulo Perforation

```

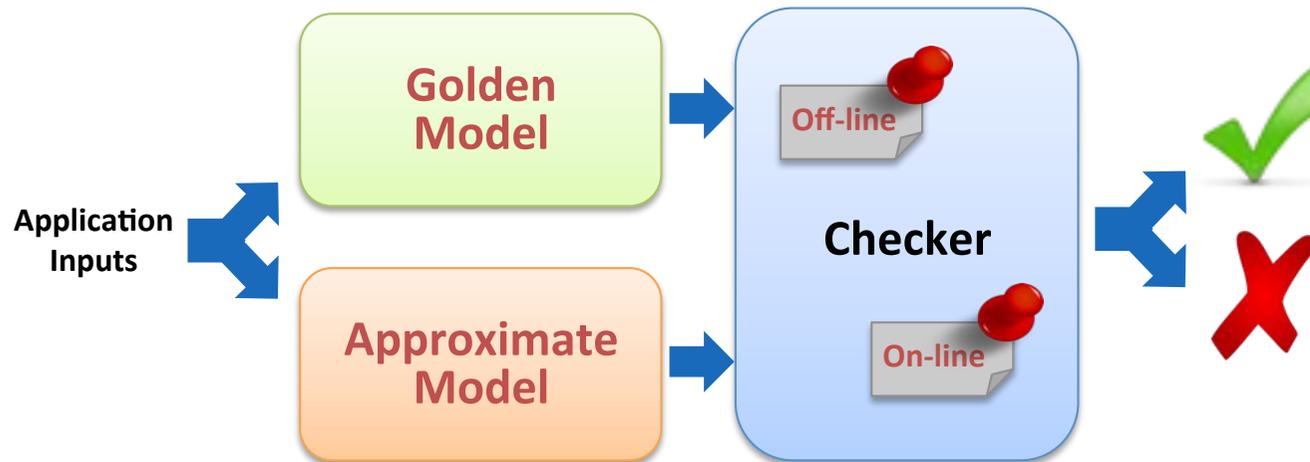
1 aspectdef FunctionCloningAndModuloPerforation
2
3 input
4   func_name,
5   modulo_value
6 end
7
8 /* The names of the new functions */
9 var approx_func_name = func_name + '_approx';
10
11 /* Clone functions */
12 select function{func_name} end
13 apply
14   exec Clone(approx_func_name);
15 end
16
17 /* Perform modulo perforation on the outermost loop */
18 select function{approx_func_name}.loop end
19 apply
20   exec ModuloPerforation(modulo_value);
21 end
22 condition
23   $loop.is_outermost
24 end
25
26 /* Change calls to the function */
27 select file.call{func_name} end
28 apply
29   var originalCode = $call.code + '!';
30   var newCode = originalCode.split(func_name).join(approx_func_name);
31
32   insert after %{
33     switch (get_best_version(/*...*/)) {
34       case 0: [[originalCode]] break;
35       case 1: [[newCode]] break;
36       default: [[originalCode]] break;
37     }
38   }%;
39 end
40
41 println('\n Modulo Perforation on the outermost loop done!');
42 end
    
```

D. Gadioli et al. "SOCRADES: A seamless online compiler and system runtime autotuning framework for energy-aware applications" DATE 2017

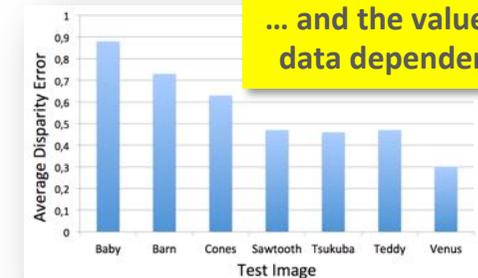
Silvano et al. "Autotuning and adaptivity in energy efficient HPC systems: the ANTAREX toolbox" Computing Frontiers 2018

Controlling the Approximation

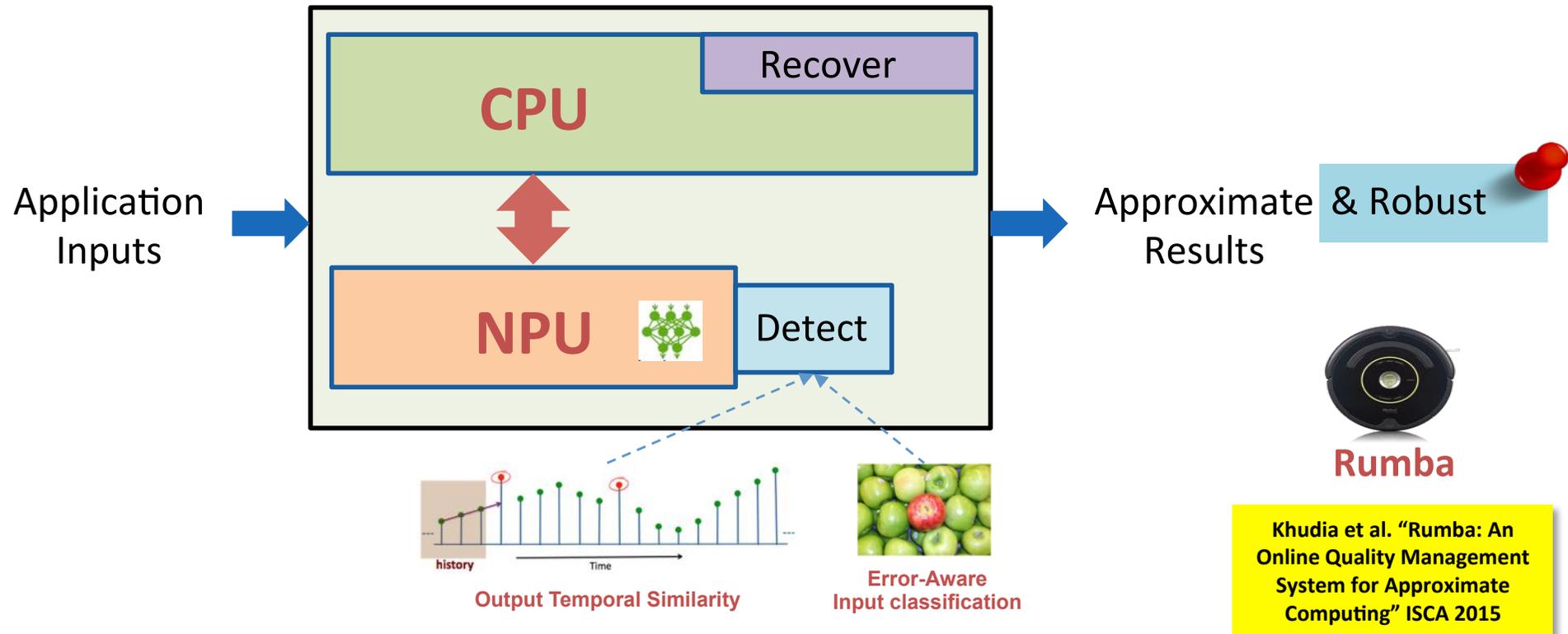
Need for monitoring the Quality of Result (QoR)



It has to be noted that establishing a suitable error measure is highly application dependent...



Online Monitoring: an Example

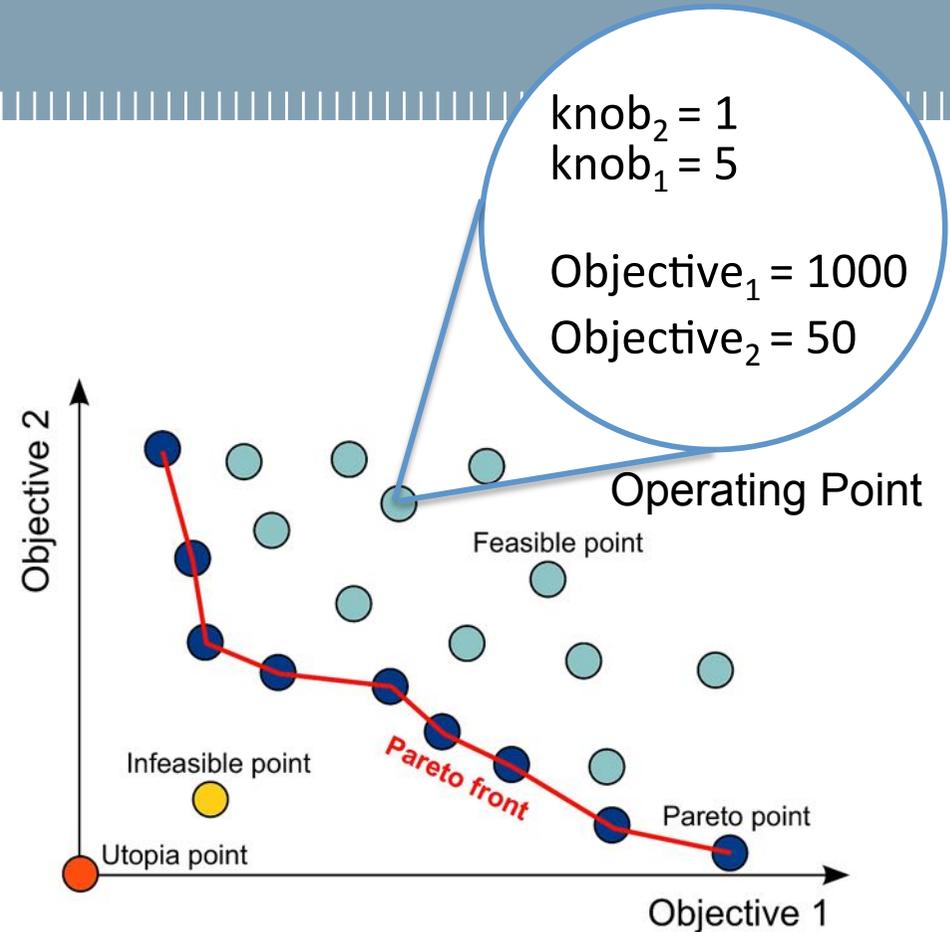


Off-line Quality Profiling

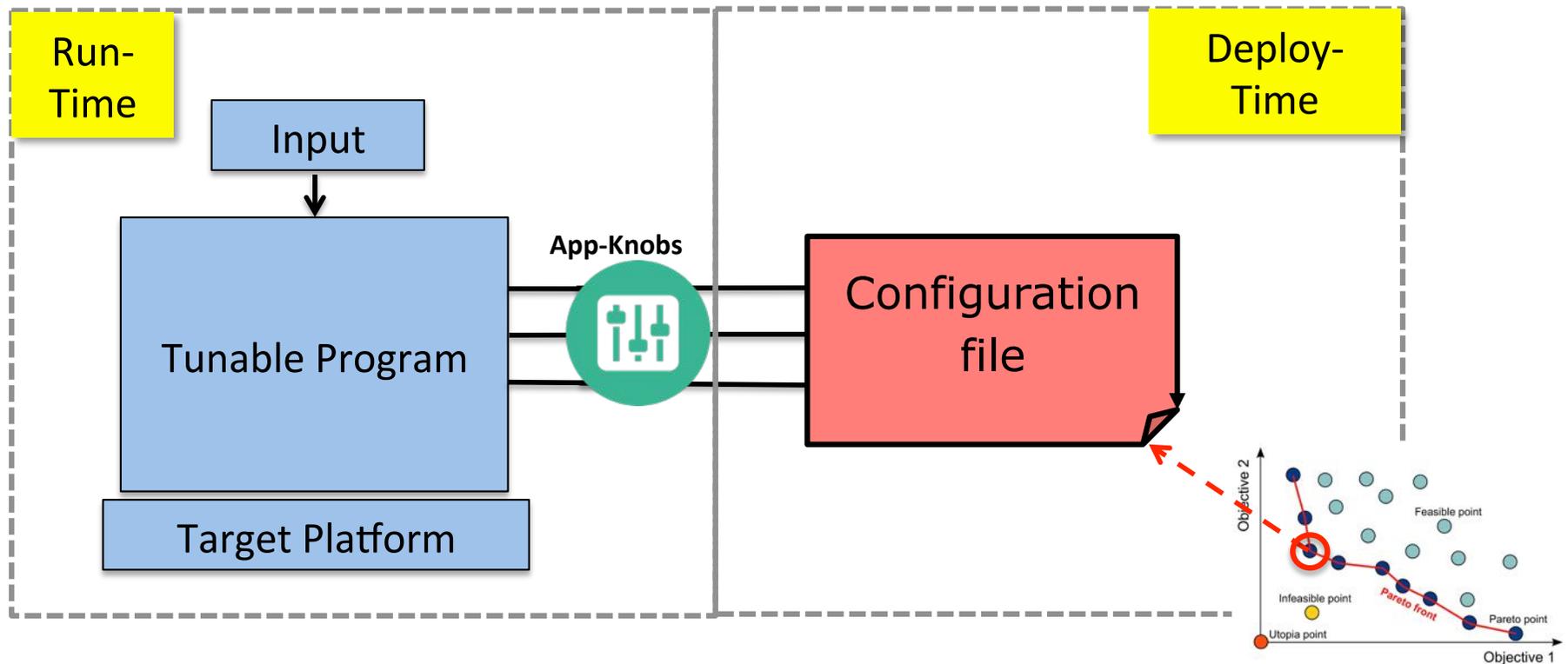
- Pure autotuning approaches
 - $\text{Error} = f(x,i)$

At **design time**:

1. **Instrument** the application
2. Perform a **Design Space Exploration**
3. **Store** the Pareto front

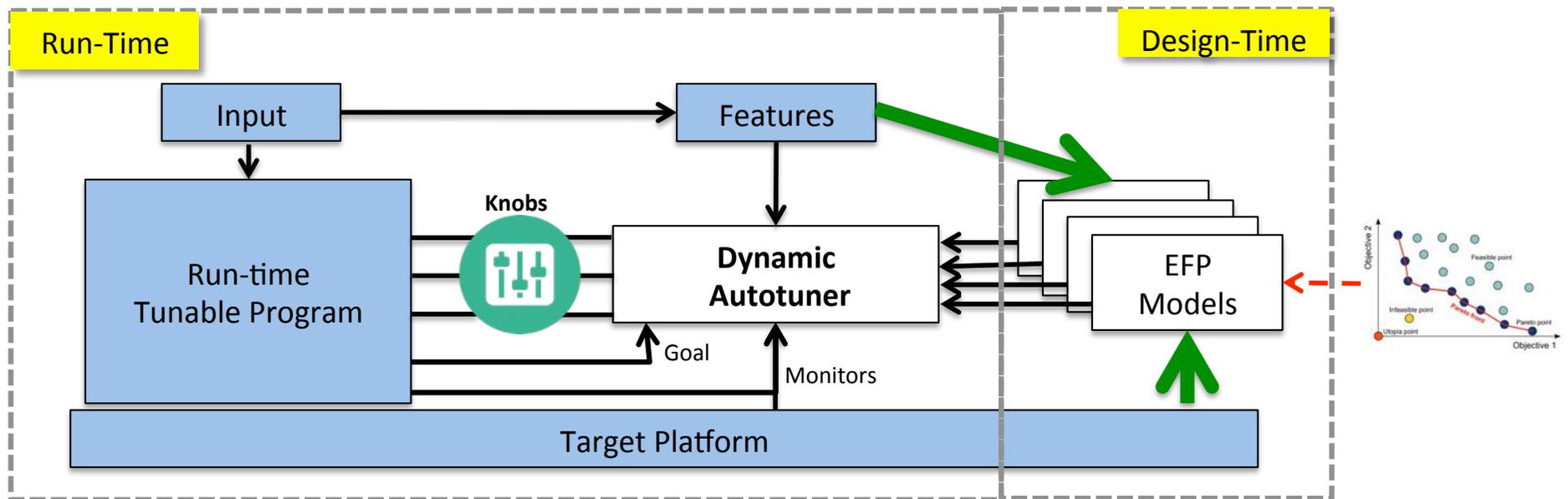


Configuring a Tunable Application



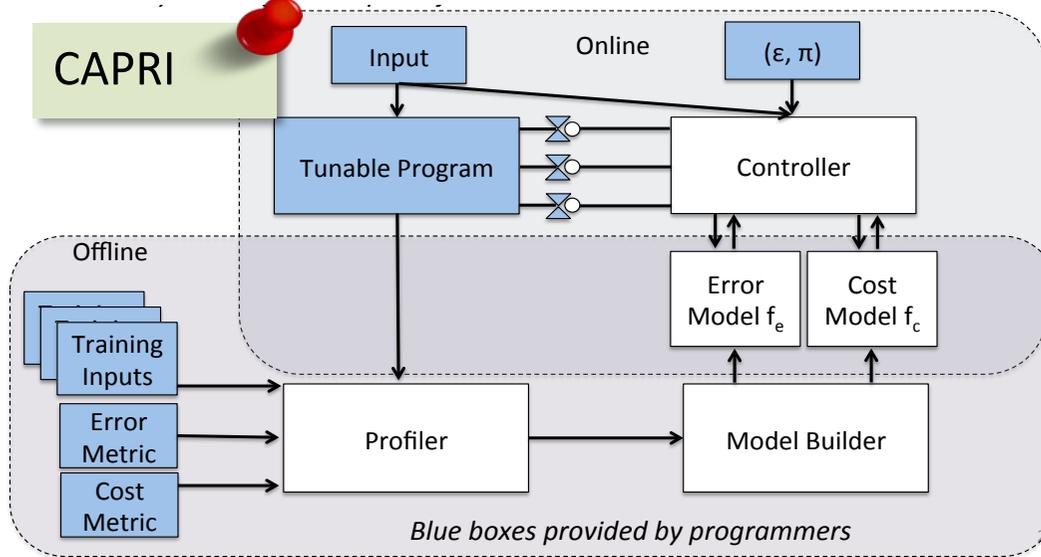
Off-line Quality Profiling (cont...)

- Pure autotuning approaches
 - $\text{Error} = f(x,i)$
- Proactive and/or reactive approaches:
 - $\text{Error} < K \Rightarrow x' : f(x',i) < K$



Off-line Quality Profiling (cont...)

- Pure autotuning approaches
 - $\text{Error} = f(x, i)$
- Proactive and/or reactive approaches:
 - $\text{Error} < K \Rightarrow x' : f(x', i) < K$

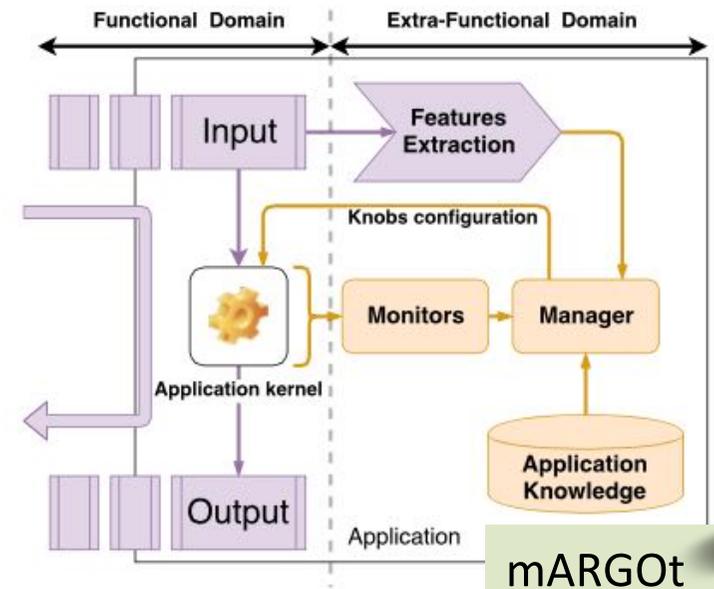
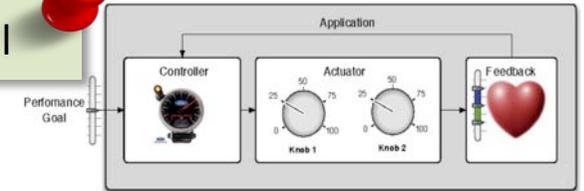


X. Sui et al. "Proactive Control of Approximate Programs" ASPLOS 2016

H. Hoffmann "JouleGuard: [...]" SOSP 2015

H. Hoffmann et al. "Dynamic knobs for responsive power-aware computing." ASPLOS 2012

PowerDial



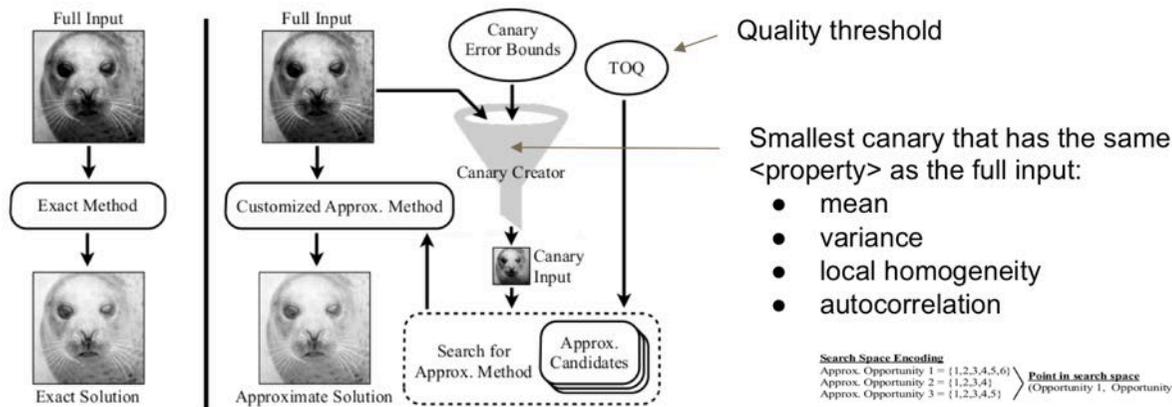
mARGOT

D. Gadioli et al. "Application autotuning to support runtime adaptivity in multicore architectures" SAMOS15
https://gitlab.com/margot_project/

Removing Offline Quality Profiling

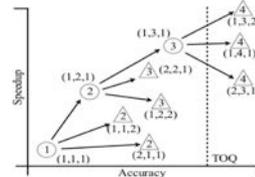
IRA – Input Responsive Approximation

It leverages **canary inputs** to select the approximation level



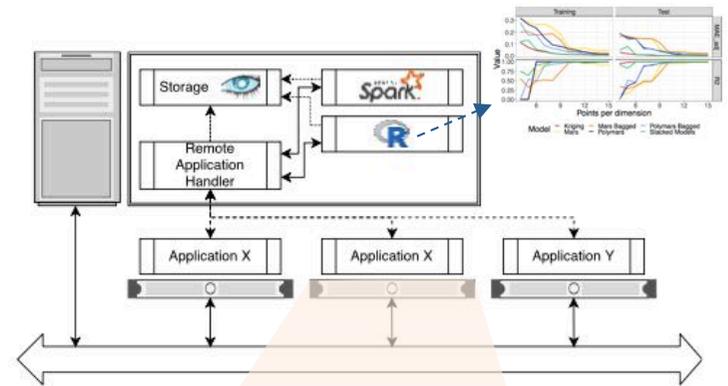
M. Laurenzano et al. "Input responsiveness: using canary inputs to dynamically steer approximation." PLDI 16.

Search Space Encoding
 Approx. Opportunity 1 = {1,2,3,4,5,6}
 Approx. Opportunity 2 = {1,2,3,4}
 Approx. Opportunity 3 = {1,2,3,4,5}

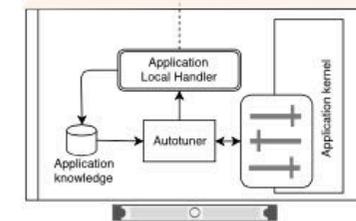


Walkthrough
Step 1 Compute (1,1,1) (1,1,1) is the baseline
Step 2 Compute (2,1,1), (1,2,1) is the steep
Step 3 Compute (2,2,1), (1,3,1) is the steep
Step 4 Compute (2,3,1), (All violate TOQ, n

mARGOt - AGORA



(a) Global structure of the distributed DSE framework



(b) Structure of an application instance

D. Gadioli et al. "mARGOt: a Dynamic Autotuning Framework Targeting Adaptivity and Controllable Approximation" SoftwareX

Conclusions... Why taking Care of Approximation?

