



Dynamic Application Autotuning for Approximate Computing

Gianluca Palermo – gianluca.palermo@polimi.it

Introduction

DYNAMIC AUTOTUNING

At Run-time
according to
different situation

Made
Automatically

Capability to customize
application knob values

It is a way to constantly improve performance/energy with low developer effort over a wide range of run-time situations

Application Tuning: Why?

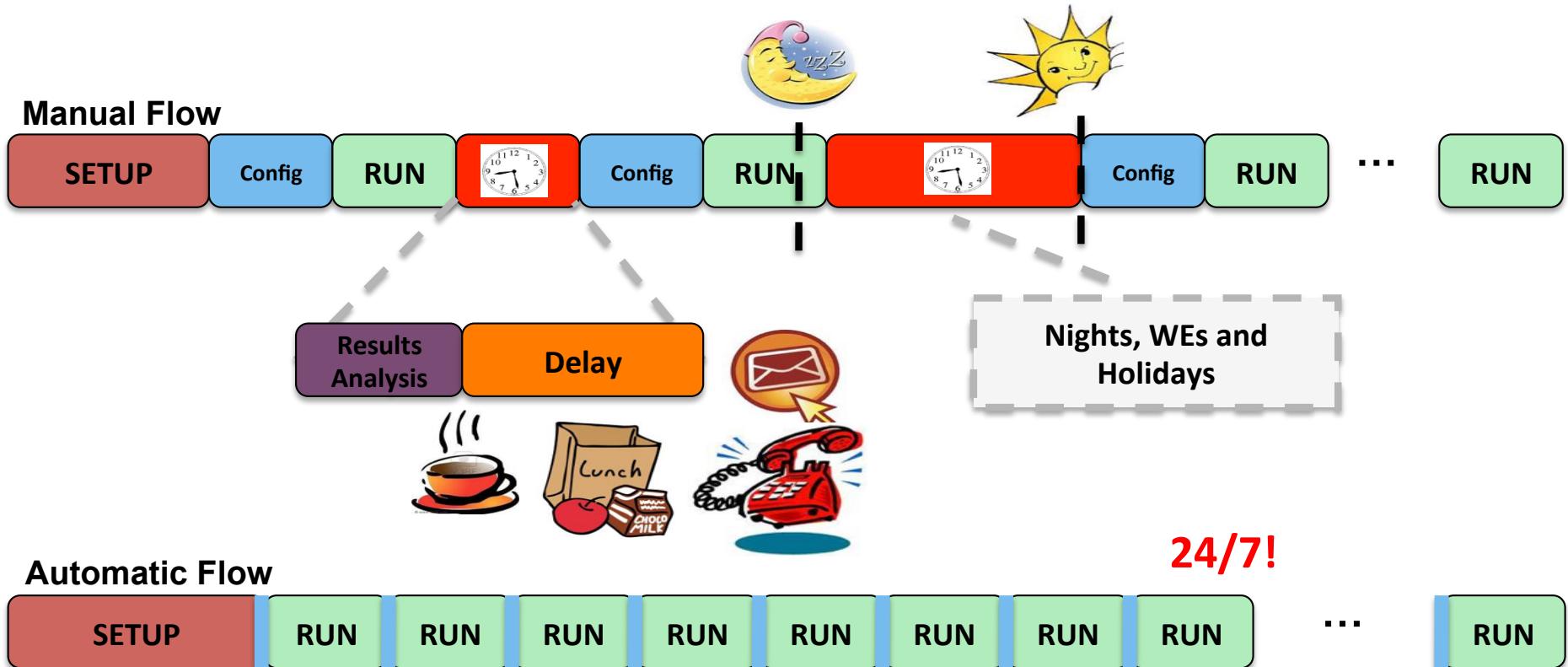
- **Applications are parameterized** to facilitate the portability and improve the productivity (*best practice*)
- Large set of application ***knobs*** with not so easy relations with the extrafunctional properties
- Complex architectures including **frequent updates and novel evolutions**

Support for performance portability on wider range of today architectures



Possibility to easily exploit upcoming (or future) platforms

Automatic Vs Manual Tuning

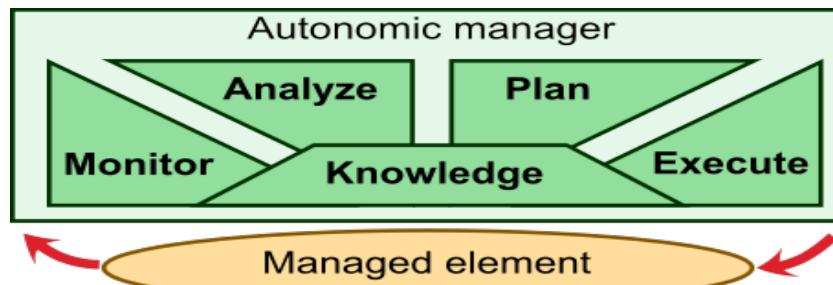


Dynamic: Why?

- The application **requirements** might **change** at runtime
 - E.g. Depending on the observed scene
- The extra-functional properties might be **input dependent**
 - Input features (e.g. input size, autocorrelation)
- The extra-functional properties might depend on the **system workload**
 - Shared computational resources
 - Core frequencies

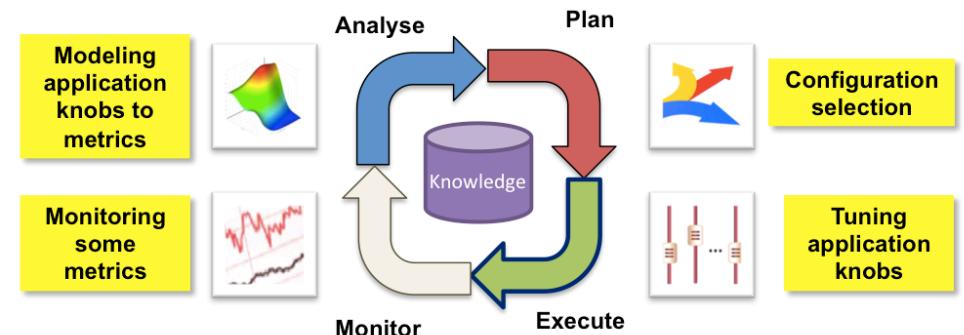
Dynamic: The autonomic computing vision

Each application is considered an **autonomic agent**, which is capable of self-management.

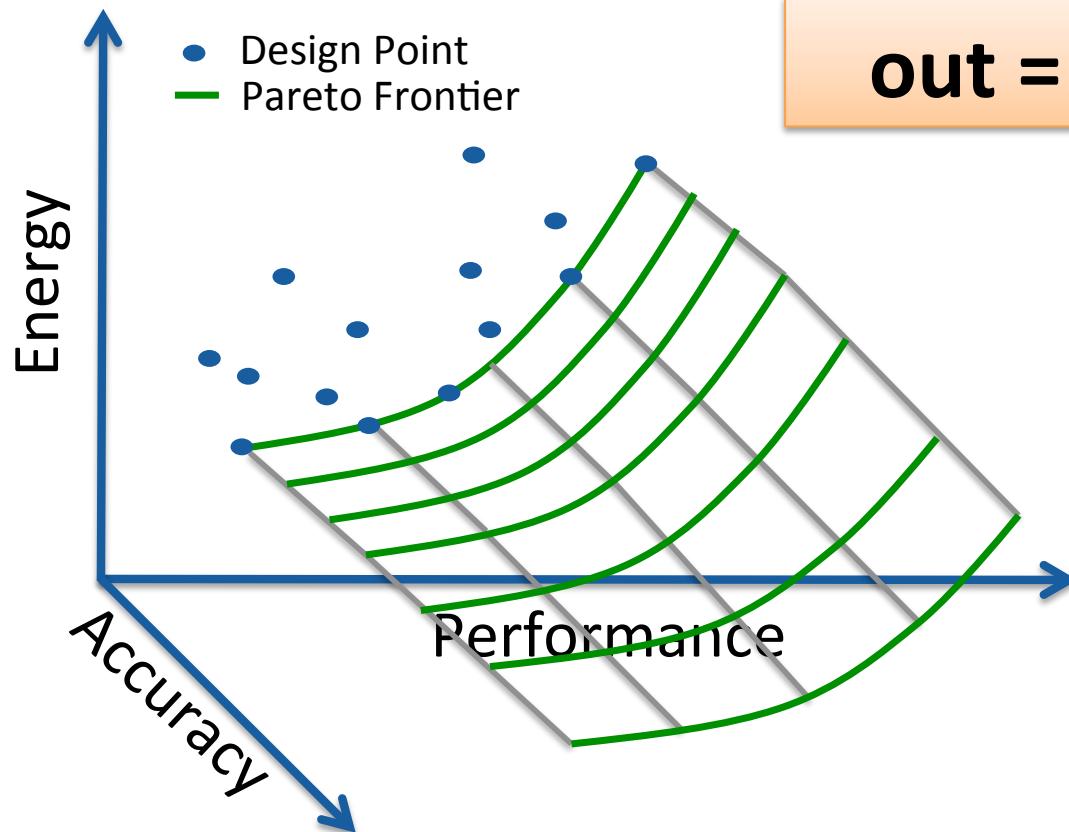


Kephart, Jeffrey O., and David M. Chess. "The vision of autonomic computing." *Computer* 36.1 (2003)

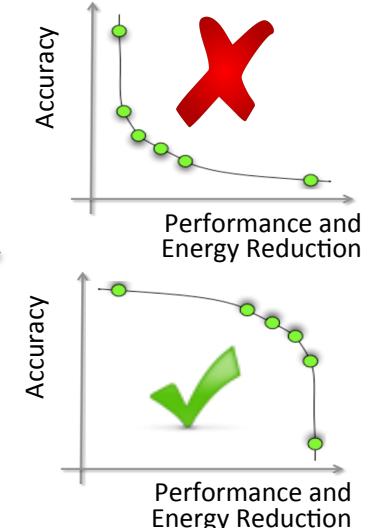
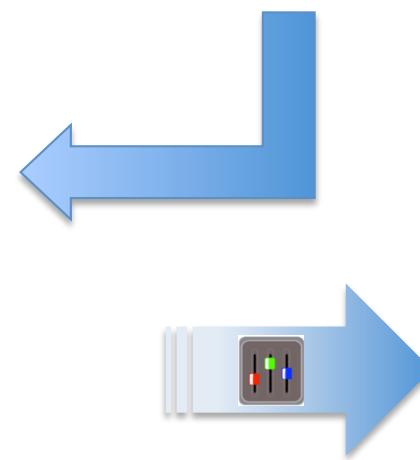
- self-configuration
- self-healing
- self-protection
- **self-optimization**



Multi-objective optimization problem



$$\text{out} = f(\text{input}, k_1, k_2, \dots, k_n)$$



Approximate Computing Applications

- Approximation adds complexity
 - Most of a code base can't 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.



Possible Application Domains...



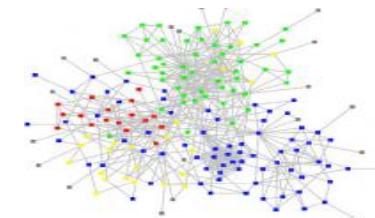
Image Processing



Robotics



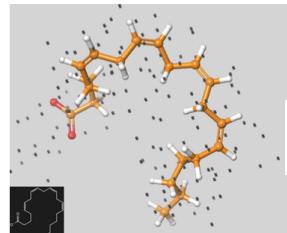
Big Data Analytics



Graph Analytics

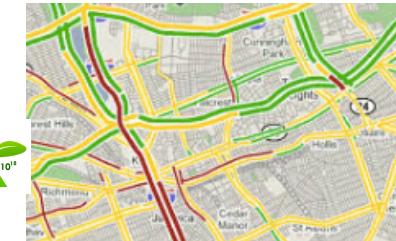


Multimedia



ANTARE^{10*}

Drug Discovery



Traffic Prediction

... where 100% of accuracy not always required

HOWEVER ...

... we have to pay attention



No Traffic



No Intrusion

Extra-Functional Properties



Functional Description



What to do...



Extra-Functional Properties



How It is done...

Extra-functional properties

$$\text{out} = f(\text{input}, k_1, k_2, \dots, k_n)$$

EFP of the output:

- Size (i.e. of the file)
- Accuracy (i.e. PSNR)
- Dimension (i.e. Resolution)

EFP of the elaboration:

- Throughput/Latency
- Power/Energy consumption
- Resource utilization

Application level Knobs

- Algorithm Selection
- Application Parameters
- Parallelism
- Code Perforation

Compiler level Knobs

- Flag selections
- Loop transformations
- Float-to-Fixed point Conversion
- Parallel Pattern Replacement

Target problem definition

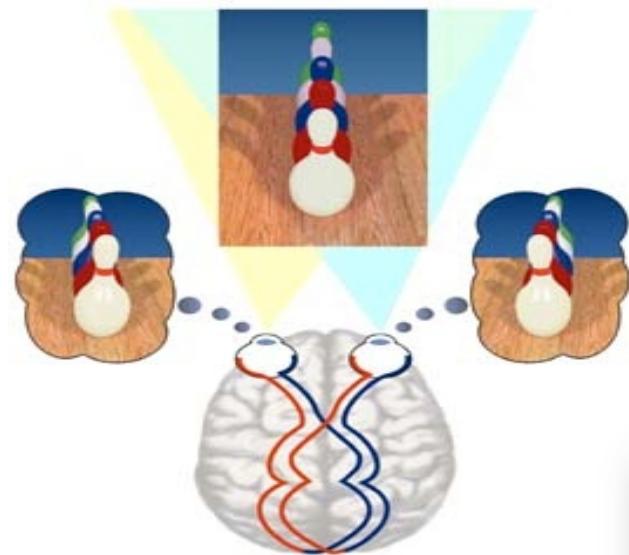
$$\text{out} = f(\text{input}, k_1, k_2, \dots, k_n)$$


Which values should we use?

The best ones? What do you mean with "best"?
One single configuration?



Let's have an example...



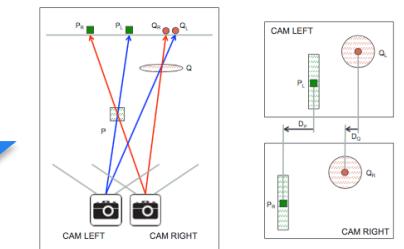
2 eyes = 3 dimensions



Left camera

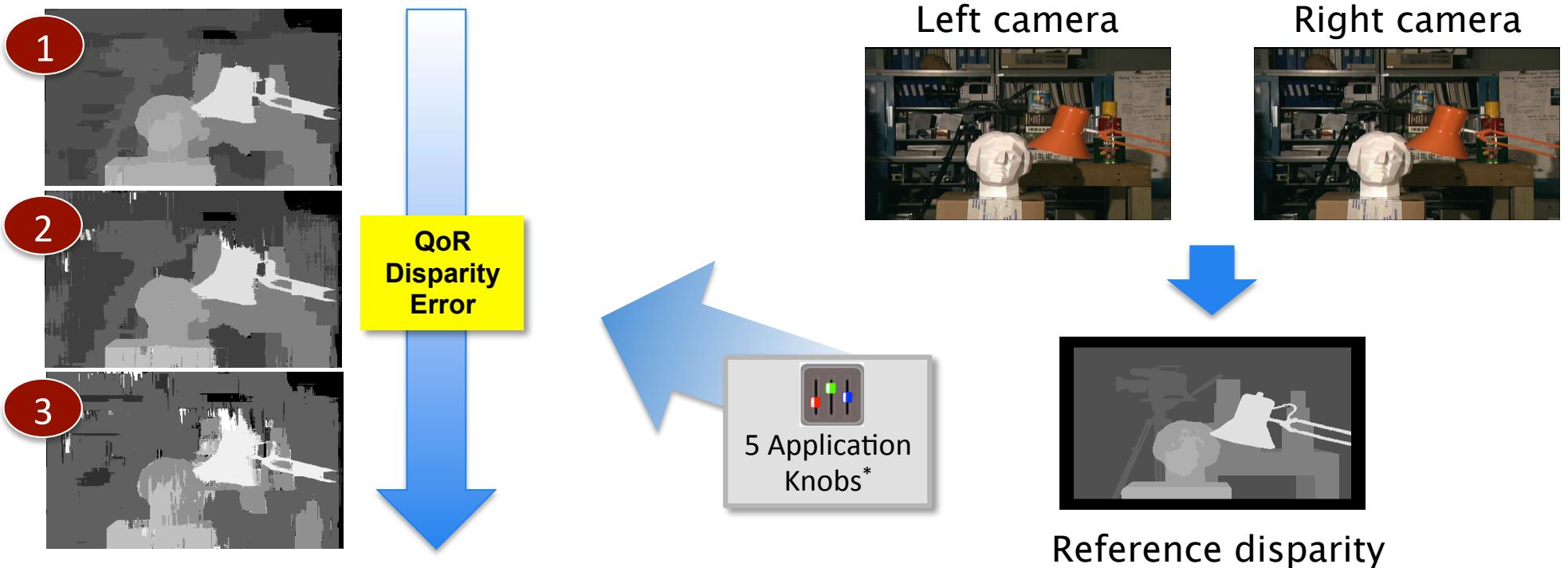


Right camera



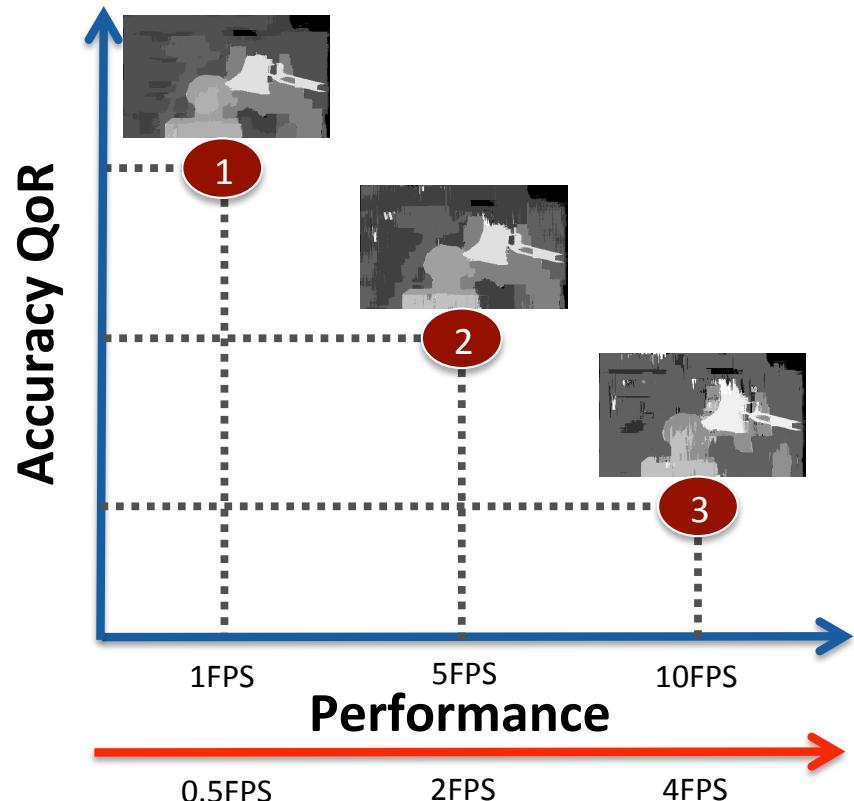
Reference disparity

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



Autotuning Problem

The definition of "best" should be expressed as a **multi-objective constrained optimization problem**

$$\text{out} = f(\text{input}, k_1, k_2, \dots, k_n)$$

Select **K1...Kn** to optimize **EFP(out/f)** subject to **EFP(out/f) constraints**, and considering **input characteristics**



POLITECNICO
MILANO 1863

Autotuning SoTa Overview

Applications Autotuning: State of the Art

Classification of autotuning approaches:

- Type of the decision
- How they are packaged into the application
- Approach to the decision
- When to apply

P. Balaprakash, J. Dongarra, T. Gamblin, M. Hall, J.K. Hollingsworth, B. Norris, R. Vuduc – Proceedings of IEEE

Approach	Definition	How Packaged	Examples
Library	A library is tuned for different architectures and input data sets.	ATLAS, OSKI, MOL, PETW, SPIRAL	
Compiler-directed	The compiler generates multiple implementations and selects among them.	CHILL, Ora, Pendrake	
Application-level	The programmer expresses parameters and code variants to select among.	ActiveHarmony, Ora, OpenTuner	
Approach to Selection			
Model-based	A performance or ML model guides selection.	Performance model, Classifier model	
Search-based	Empirical search to identify selected.	Heuristic search	
Hybrid	A model prunes search space, followed by search	Prune undesirable areas of search space	
Types of Decisions			
Algorithm selection	Select among alternative algorithms that are functionally equivalent.	PEINC, Trilinos, LightHouse	
Code variant selection	Select among fundamentally different implementations of an algorithm.	Data layout/representation	
Parallelization	Select among different parallelization strategies.	SIMD/x86, OpenMP vs. CUDA, multi-level	
Code transformation	Select among different code transformation sequences.	tile (tiled), Fusion (tiled), Wavefront (warp)	
Parameter tuning	Adjust context-specific parameters.	# of tasks/threads, tile size, tiling factor	
When to Apply			
During coding	Library/application migrates to new architecture.	Adjust data layout and parameters	
Offline at runtime	Code is generated or modified explicitly.	Code generation or recompilation	
Online, incrementally	Automatic decisions implemented over application runs.	Expert search space for context	
Dynamic, runtime	Context model to select appropriate implementation.	Based on input features	
Integration into Application			
Tuned code	Tuned code is inserted into the application.	Architecture-specific code	
Selectable code	Multiple implementations and a selector.	Input-dependent selection	
Dynamic	Code is dynamically generated.	Too many variants	
Full tuning integration	Tuning occurs as part of build.	Compact code and forward scalability	
Runtime Measurement			
Instrumentation	Gather standard measurements transparently to the user.	Dyninst, TAU, HPC Toolkit	
Semantic annotation	Tools or user direct instrumentation.	Caliper	
Performance database	Collect measurements in a database.	TAUdb	

Type of decisions

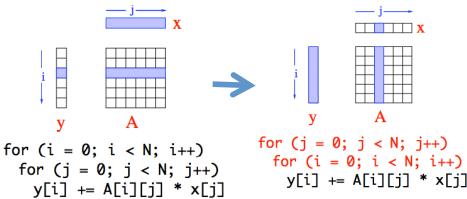
- Compiler options
- Code transformations and data layout:
 - i.e. Fusion, Interchange, Unrolling, Tiling



A. Ashouri, W. Killian, J. Cavazos, G. Palermo, and C. Silvano. A Survey on Compiler Autotuning using Machine Learning. ACM Computing Surveys (CSUR)

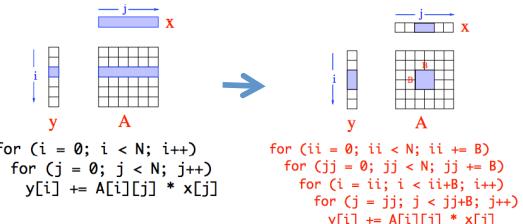
```

do i = 1, n           do i = 1, n
    c[i] = a[i]       c[i] = a[i]
end do                 end do
do i = 1, n           b[i] = a[i]
    b[i] = a[i]       end do
end do
for (i = 0; i < N; i++)
    for (j = 0; j < N; j++)
        y[i] += A[i][j] * x[j]
    
```



```

for (i = 0; i < N; i++)
    A[i] = ...
    ↓
    for (i = 0; i < N; i += 4)
        A[i] = ...
        A[i+1] = ...
        A[i+2] = ...
        A[i+3] = ...
        for (j = 0; j < N; j++)
            for (i = 0; i < N; i++)
                y[i] += A[i][j] * x[j]
    
```

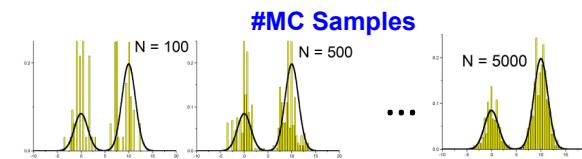
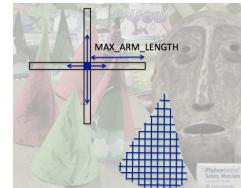


- Algorithm selection
 - Select among different implementations of functionally equivalent versions of the algorithm (i.e. for sorting: bubble, merge, quick, heap, radix, bucket...)
- Parallelization:
 - Parallel patterns, #workers, OmpThread, multi-level parallelism parameters
 - ...Approximation parameters

J. Ansel et al. "PETABRICKS: a language and compiler for algorithmic choice" PLDI09

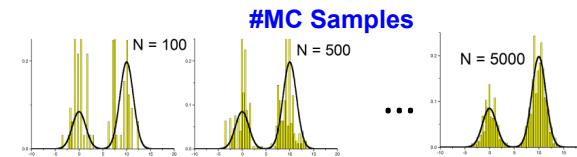
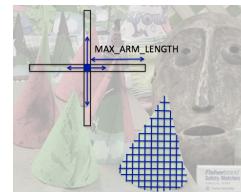
Type of decisions: Application-level approximations

- *Application parameters*

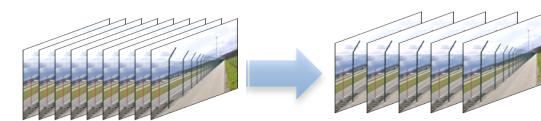


Type of decisions: Application-level approximations

- *Application parameters*

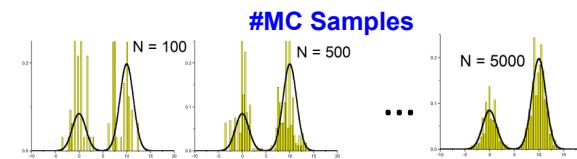
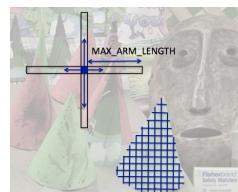


- *Task skipping*

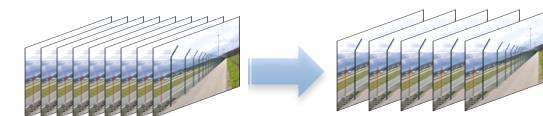


Type of decisions: Application-level approximations

- *Application parameters*

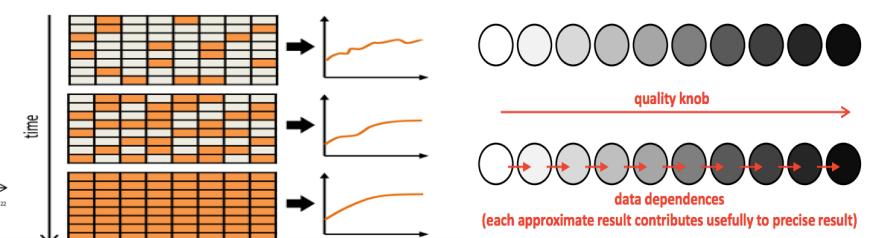
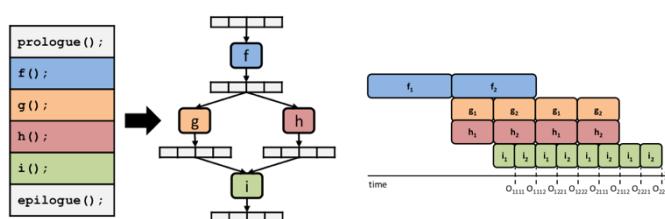


- *Task skipping*



Computation model that represents an approximate application as a pipeline

M.J. San et al "The anytime automaton." ACM SIGARCH Computer Architecture News. 2016.

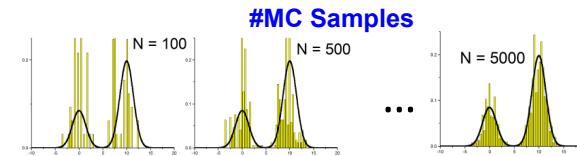


V. Vassiliadis et al "Exploiting Significance of Computations for Energy-Constrained Approximate Computing" IJPP 2016.

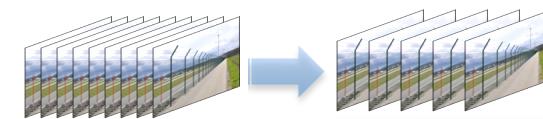
```
#pragma omp task [significant(...)] [label(...)]
[in(...)] [out(...)] [approxfun(function())]
```

Type of decisions: Application-level approximations

- *Application parameters*



- *Task skipping*



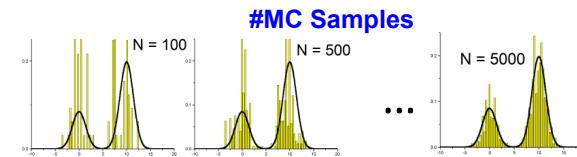
```
switch (get_quality_level(/*...*/)) {  
    case 0: foo(...); break;  
    case 1: foo_approx_Q1(...); break;  
    case 2: foo_approx_Q2(...); break;  
    /* ... */
```

- *Multiversioning*

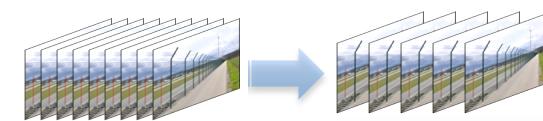
- Considering different performance/accuracy trade-offs

Type of decisions: Application-level approximations

- *Application parameters*



- *Task skipping*



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

- *Multiversioning*

- Considering different performance/accuracy trade-offs

DSL or annotation based approaches

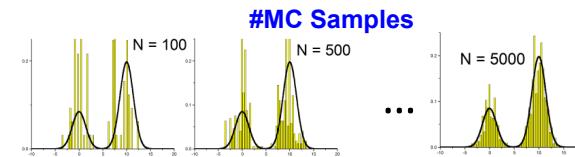
J Ansel et al. "PetaBricks: a language and compiler for algorithmic choice" PLDI 2009

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

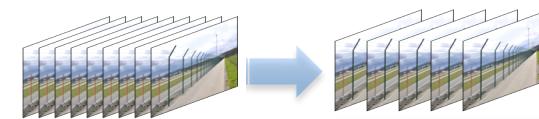
```
1 transform Sort  
2 from A[n]  
3 to B[n]  
4 {  
5     from(A a) to(B b) {  
6         tunable WAYS;  
7         /* Mergesort */  
8     } or {  
9         /* Insertionsort */  
10    } or {  
11        /* Radixsort */  
12    } or {  
13        /* Quicksort */  
14    }  
15 }
```

Type of decisions: Application-level approximations

- *Application parameters*



- *Task skipping*



```
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;
}
```

- *Multiversioning*

- Considering different performance/accuracy trade-offs

- *Approximate Memoization*

- User partial key for the lookup

($p \& 0xffff0000$)

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

Type of decisions: Compiler approaches

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

Type of decisions: Compiler approaches

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

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

The diagram illustrates the transformation of floating-point operations in C code into their equivalents using the MPFR library. It shows two columns: 'Original code' and 'Rewritten code'. In the 'Original code' column, lines 2 and 3 are highlighted with blue boxes, indicating they are being converted. In the 'Rewritten code' column, the corresponding MPFR library calls are shown, with the precision value '53' circled in orange. A callout box labeled 'Precision assigned at runtime' points to this circled value. Arrows from the highlighted lines in the original code point to the corresponding MPFR calls in the rewritten code.

```
Original code:
1. function(){
2.     double var1 = 1.0;
3.     double var2 = var1 + 1.0;
4.     ...
5. }
```

```
Rewritten code:
1. function(){
2.     mpfr_t var1;
3.     mpfr_init2(var1, 53);
4.     mpfr_set_d(var1, 1.0, MPFR_RNDN);
5.     mpfr_t var2;
6.     mpfr_init2(var2, 53);
7.     mpfr_add_d(var2, var1, 1.0, MPFR_RNDN);
8.     ...
9. }
```

Type of decisions: 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.

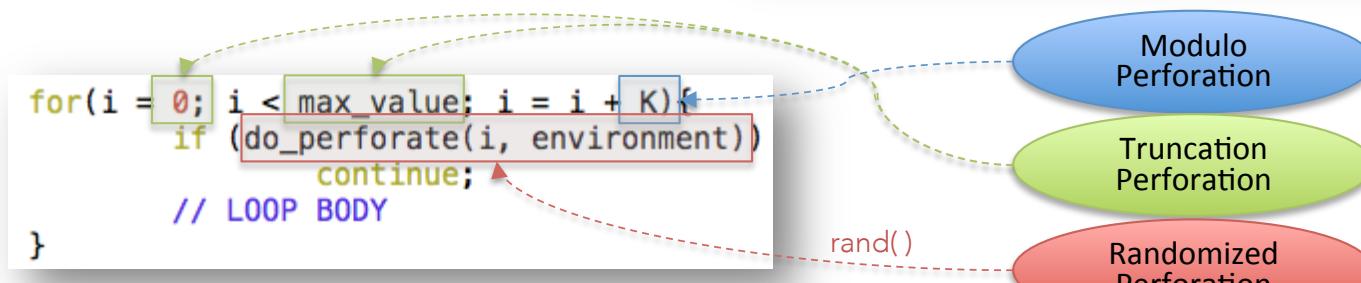
The diagram illustrates the transformation of floating-point code. On the left, the 'Original code' is shown as:

```
1. function(){
2.     double var1 = 1.0;
3.     double var2 = var1 + 1.0;
4.     ...
5. }
```

On the right, the 'Rewritten code' is shown using the MPFR library:

```
1. function(){
2.     mpfr_t var1;
3.     mpfr_init2(var1, 53);
4.     mpfr_set_d(var1, 1.0, MPFR_RNDN);
5.     mpfr_t var2;
6.     mpfr_init2(var2, 53);
7.     mpfr_add_d(var2, var1, 1.0, MPFR_RNDN);
8.     ...
9. }
```

A callout box labeled 'Precision assigned at runtime' points to the precision parameters (53) in the rewritten code.



How they are packaged

- Library-based
 - Mainly Performance-critical routines
 - Parameterization and tuning engine are not exposed
- Compiler-Level
 - The code transformations and variants are injected by the compiler/code generation framework
 - The tuning engine can be embedded into the compiler or parameters are exposed for tuning
- Application-Level
 - The programmer parameterizes the application
 - External tuning frameworks

R. Whaley, J. Dongarra, "ATLAS: Automatically tuned linear algebra software," SC98

M. Frigo, S. Johnson, "The design and implementation of FFTW3," Proc. of IEEE 2005

R. Vuduc, et al. "OSKI: A library of automatically tuned sparse matrix kernels," Journal. of Physics, 2005

M. Puschel, et al. "SPIRAL: A generator for platform-adapted libraries of signal processing algorithms," JHPCA04

A. Tiwari et al. "CHILL: scalable auto-tuning framework for compiler optimization," IPDPS09

B. Norris, et al. "Annotation-based empirical performance tuning using ORIO," IPDPS09

J. Ansel et al. "PETABRICKS: a language and compiler for algorithmic choice" PLDI09

D. Gadioli, et al. "SOCRATES - A seamless online compiler and system runtime autotuning framework" DATE18

J. Ansel et al. "OPENTUNER: An extensible framework for program autotuning," PACT14

D. Gadioli et al. "mARGOT: a Dynamic Autotuning Framework Targeting Adaptivity and Controllable Approximation" TCOM19

Approach to the decision

- Search-Based
 - Complete enumeration (small configuration space)
 - Heuristics used to navigate the configuration space
 - Genetic algorithms, simulated annealing, particle swarm optimization
 - Mix of global and local searches
- Model-based
 - Pure analytic model of the performance
 - Requires a huge knowledge of the target problem
 - Source code based model
 - Based on static code analysis
 - Requires to generate the modified code
 - Training-based
 - Subset of configurations are evaluated and a predictive model is built using Machine Learning techniques
- Hybrid
 - Sequence of Model- and Search-based techniques

A. Rasch, et al. "ATF: A Generic Auto-Tuning Framework," HPCC17

R. Miceli et al., "AUTOTUNE A plugin-driven approach to the automatic tuning of parallel applications" 2012

B. Norris, et al. "... ORIO," IPDPS09

R. Whaley, J. Dongarra, "ATLAS..." SC98

J. Ansel et al."OPENTUNER ... , " PACT14

G. Mariani, et al. "Scaling Properties of Parallel Applications to Exascale." IJPP16

P. Balaprakash et al. "Can search algorithms save large-scale automatic performance tuning?" ICCS11

J. Bergstra, et al. "Machine learning for predictive auto-tuning with boosted regression trees," INPAR12

T. Martinovic, et al. "On-line Application Autotuning Exploiting Ensemble Models" arXiv19

G. Palermo, et al. "ReSPIR: A Response Surface-based Pareto Iterative Refinement for DSE". TCAD09

G. Mariani, et al. "OSCAR: An Optimization Methodology Exploiting Spatial Correlation" TCAD12

M. Zuluaga,"SMART design space sampling to predict Pareto-optimal solutions." LCTES12

When to apply

- During Porting
 - Mainly for libraries
 - The code is tuned to adjust data layout and parameters to the new architecture
- Offline
 - The configuration space is profiled and the optimal configuration is applied
- Dynamically at RunTime
 - Autotuning decision can be improved during the application runs
 - Possibility to use a larger production environment
 - Decisions are taken considering input data characteristics and execution context

R. Whaley, J. Dongarra, "ATLAS..." SC98

M. Frigo, et al. "... FFTW3," Proc. of IEEE 2005

R. Vuduc, et al. "OSKI: ..." Journal. of Physics, 2005

M. Puschel, et al. "SPIRAL: ..." JHPCA04

J. Ansel et al. "OPENTUNER: An extensible framework for program autotuning," PACT14

B. Norris, et al. "Annotation-based empirical performance tuning using ORIO," IPDPS09

J. Ansel et al. "PETABRICKS: a language and compiler for algorithmic choice" PLDI09

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

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

D. Gadioli et al. "mARGOT: a Dynamic Autotuning Framework ..." TCOM19



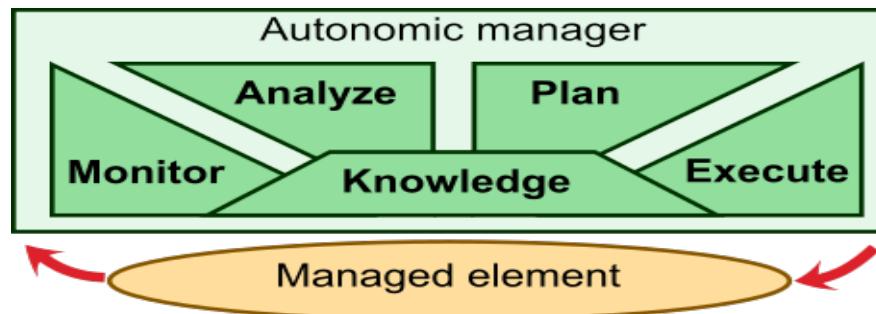
POLITECNICO
MILANO 1863

Off-line Profiling and On-Line Monitoring

The Knowledge is the key

The Knowledge of the system behaviour is the key point for any decision

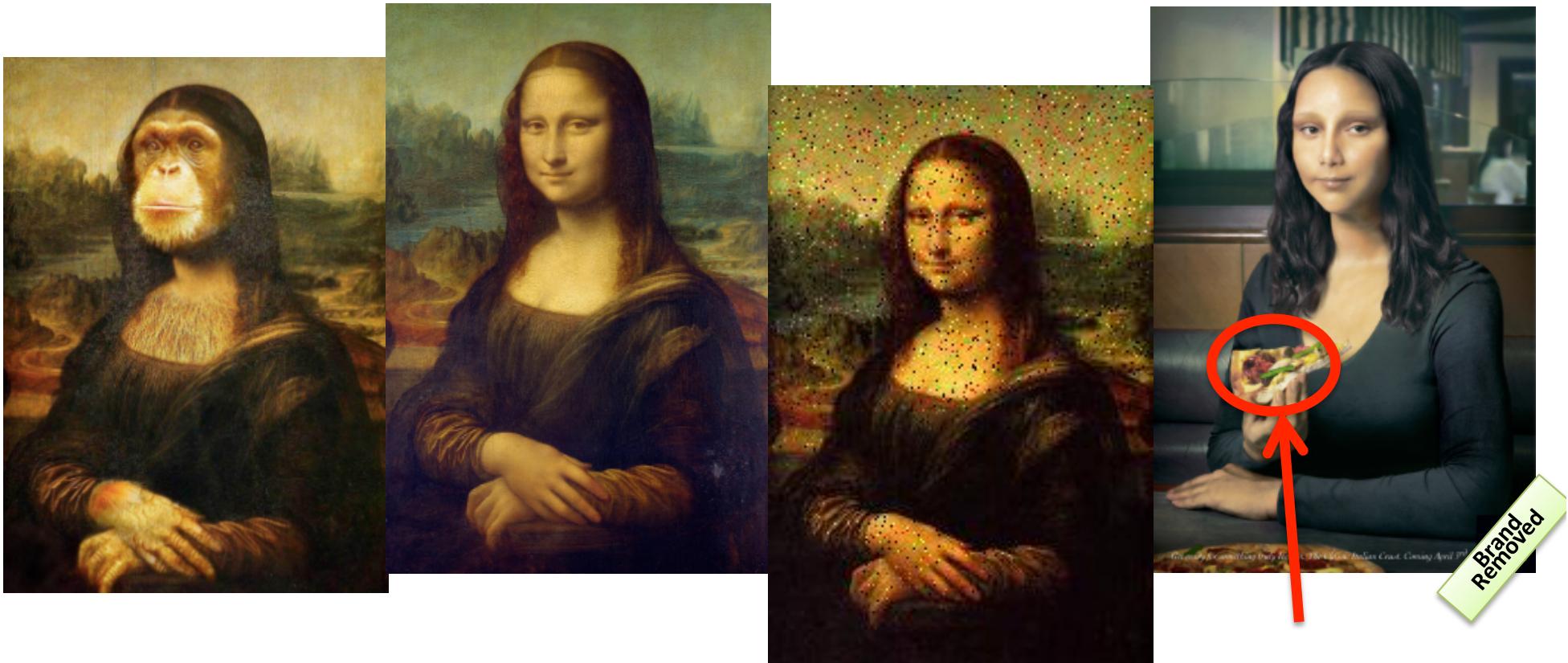
- Offline case: How to configure?
- Online case: How to re-configure?



Energy/Power/
Latency are easily
measurable...

What about
approximation
effects?!

Why taking care of approximation effects?



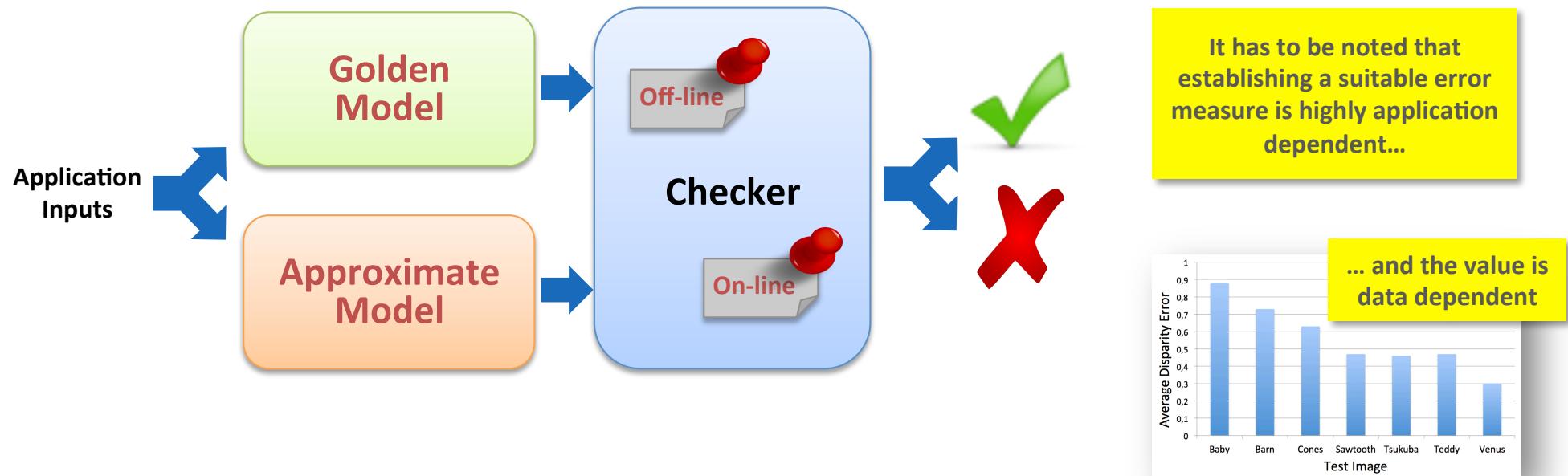
Controlling the Approximation

Need for monitoring/profiling the Quality of Result (QoR)

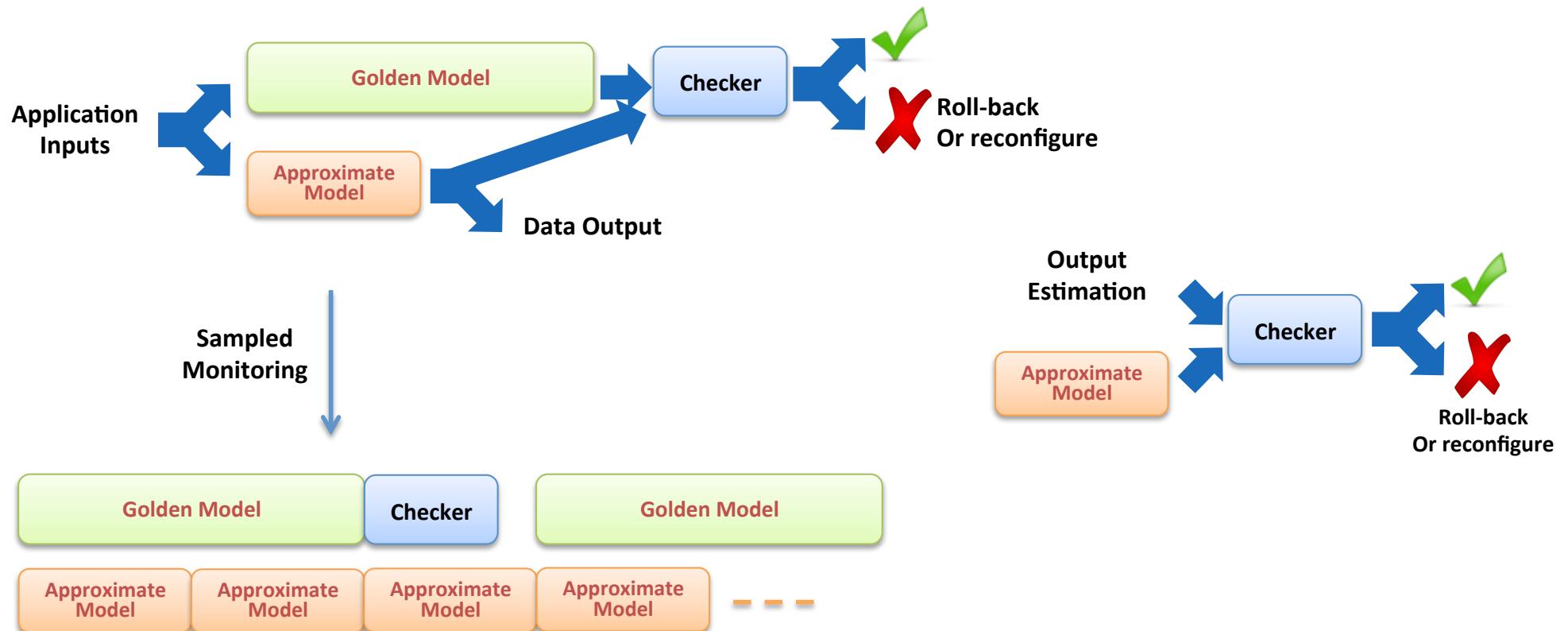


Controlling the Approximation

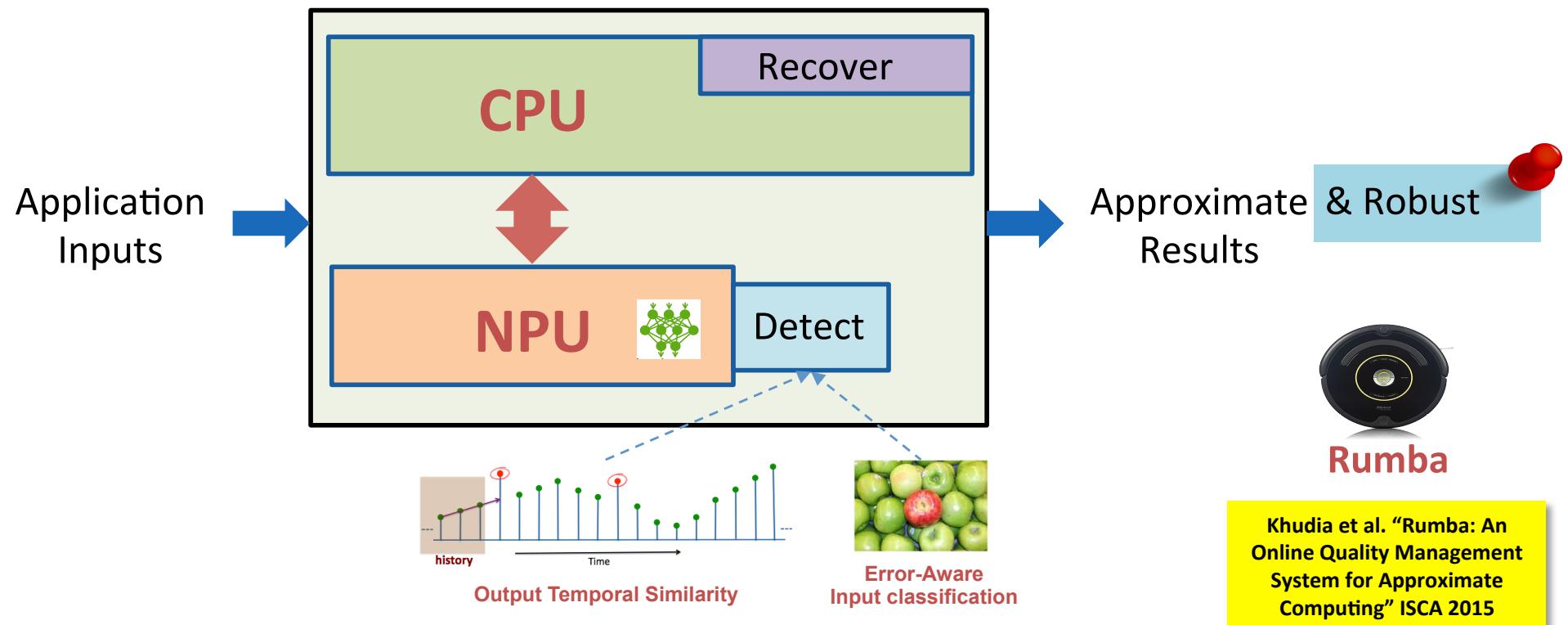
Need for monitoring/profiling the Quality of Result (QoR)



OnLine Monitoring



Online Monitoring: an example



Off-line Quality Profiling

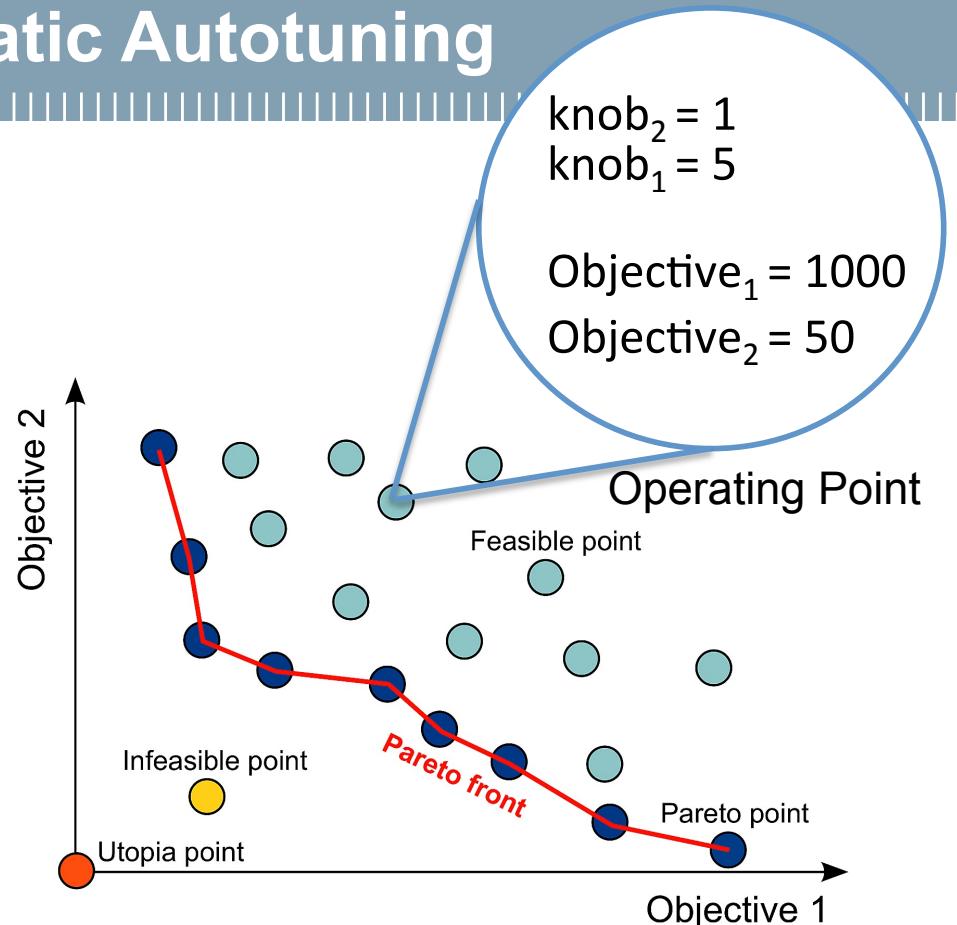
- Pure autotuning approaches
 - Error = $f(x,i)$

Off-line Quality Profiling – Static Autotuning

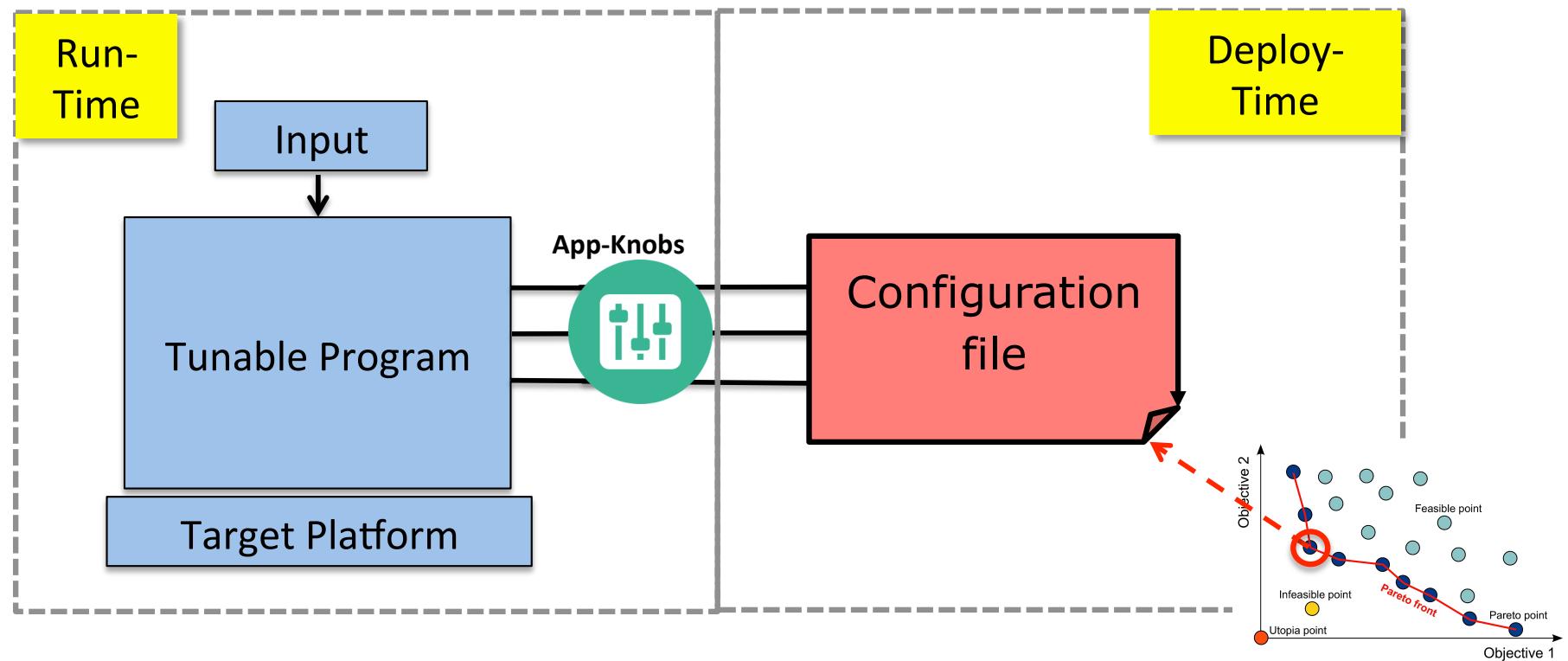
- Pure autotuning approaches
 - Error = $f(x,i)$

At **design time**:

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

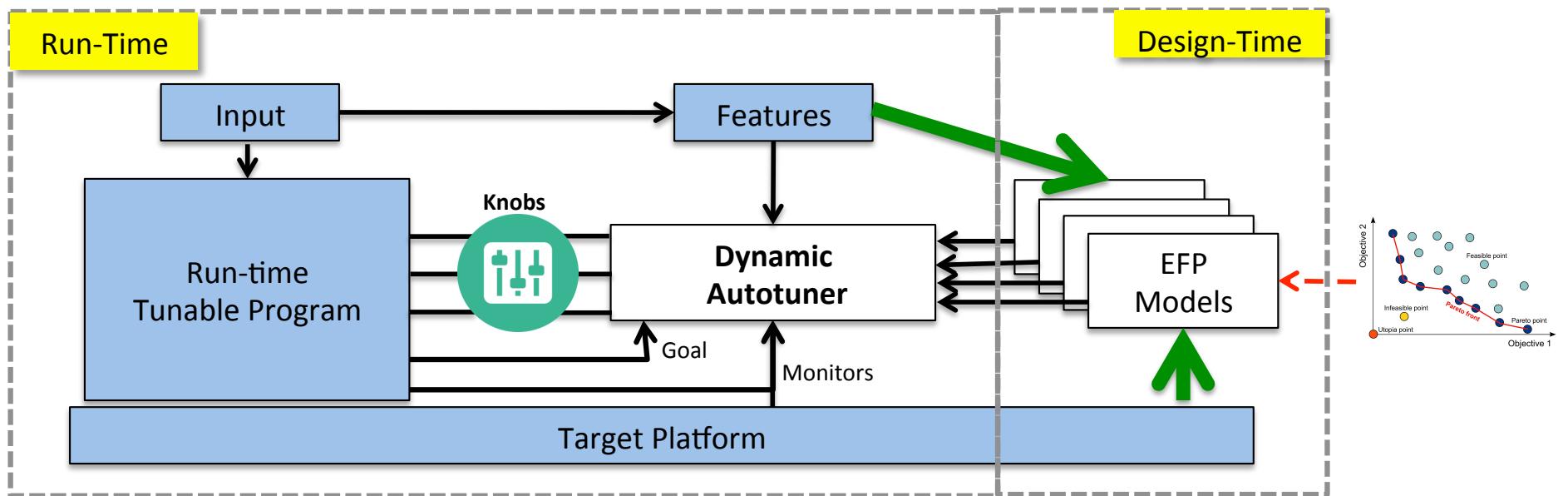


Configuring a Tunable Application



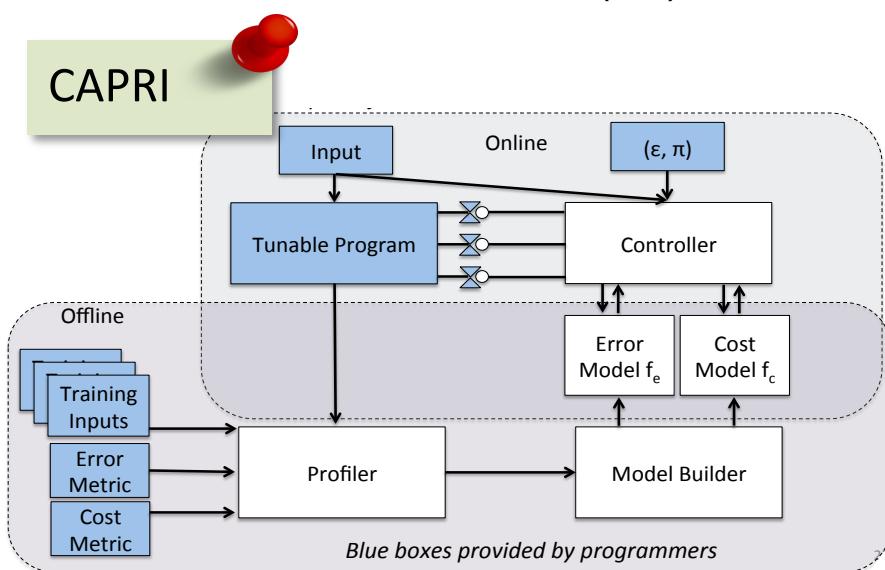
Off-line Quality Profiling – Dynamic Autotuning

- 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 – Dynamic Autotuning

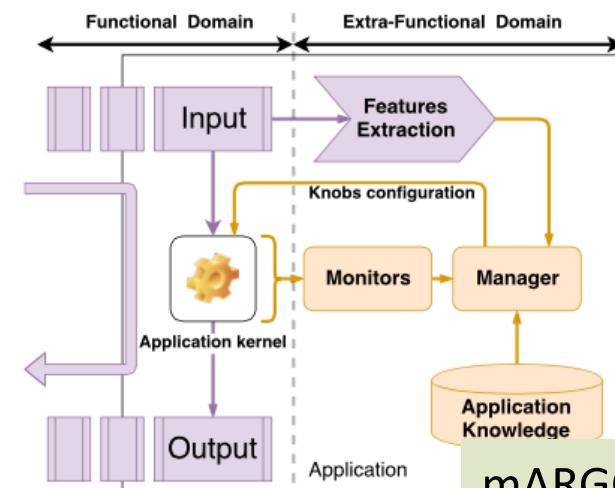
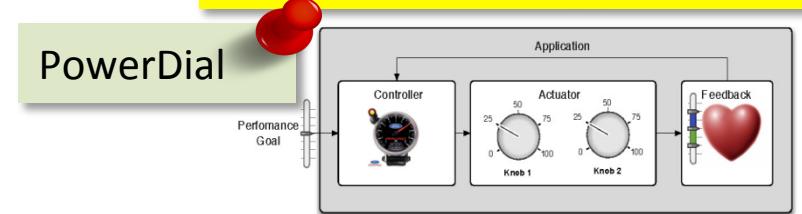
- Pure autotuning approaches
 - Error = $f(x,i)$
- Proactive and/or reactive approaches:
 - Error < K => $x' : f(x',i) < K$



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

44

H. Hoffmann et al. "Dynamic knobs for responsive power-aware computing." ASPLOS 2012; H. Hoffmann "JouleGuard: [...]" SOSP 2015

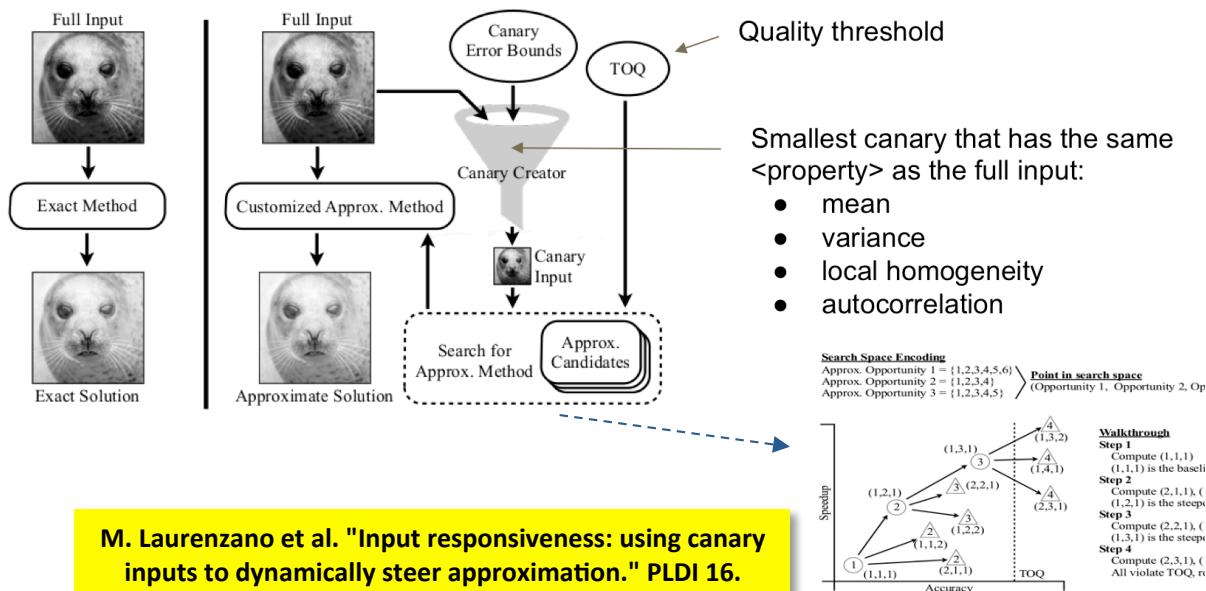


D. Gadioli et al. "mARGOT: a Dynamic Autotuning Framework Targeting Adaptivity and Controllable Approximation" TCOM19

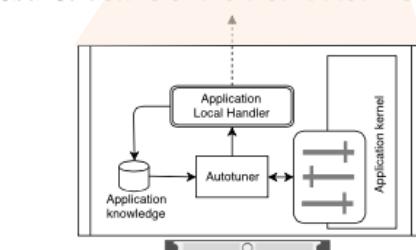
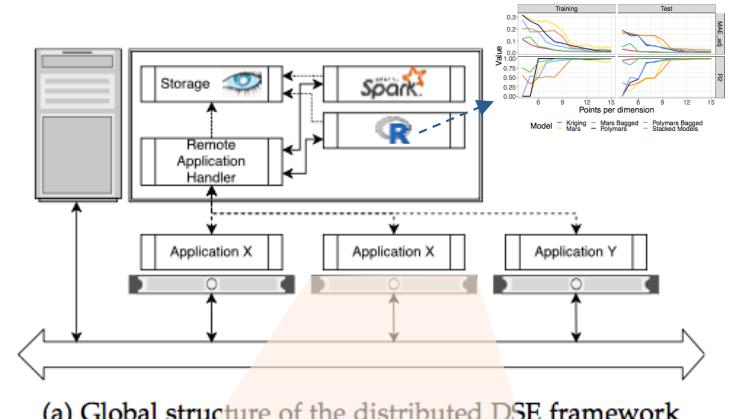
Removing Offline Quality Profiling

IRA – Input Responsive Approximation

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



mARGOT - AGORA



T. Martinovic et al. "On-line Application Autotuning Exploiting Ensemble Models" arXiv:1901.06228

Thanks also to...



ANTAREX^{10¹⁸}



Horizon 2020
European Union funding
for Research & Innovation

<http://www.antarex-project.eu/>