Building Adaptable **Software for Resilience**

Milind Kulkarni CRISP Seminar, 9/27/2018



CRISP

what do we mean by resilience?

- When people think software resilience, they often think of specific issues:
 - Fault tolerance: The ability to tolerate failures during application execution (hardware failure, failure of other software components, etc.)
 - Maybe the most common sense. Microsoft defines resiliency as "the ability of a system to gracefully handle and recover from failures"
 - Security: The ability to withstand malicious users or environment





a broader view

resilience, noun: an ability to recover from or easily adjust to misfortune or change





a broader view

resilience, noun: an ability to recover from or easily adjust to misfortune or change

software resilience, noun: an ability for software to adapt to change





software resilience as adaptability

- Software should be able to adapt to different circumstances
 - Hardware platform not the same as expected: *failing* hardware, removing hardware, adding hardware, changing hardware
 - Usage scenario not the same as expected: input is *erroneous*, input is *malicious*, input is *larger* or *smaller* than expected

• Goal: software written one time, but can adapt to a wide variety of different scenarios





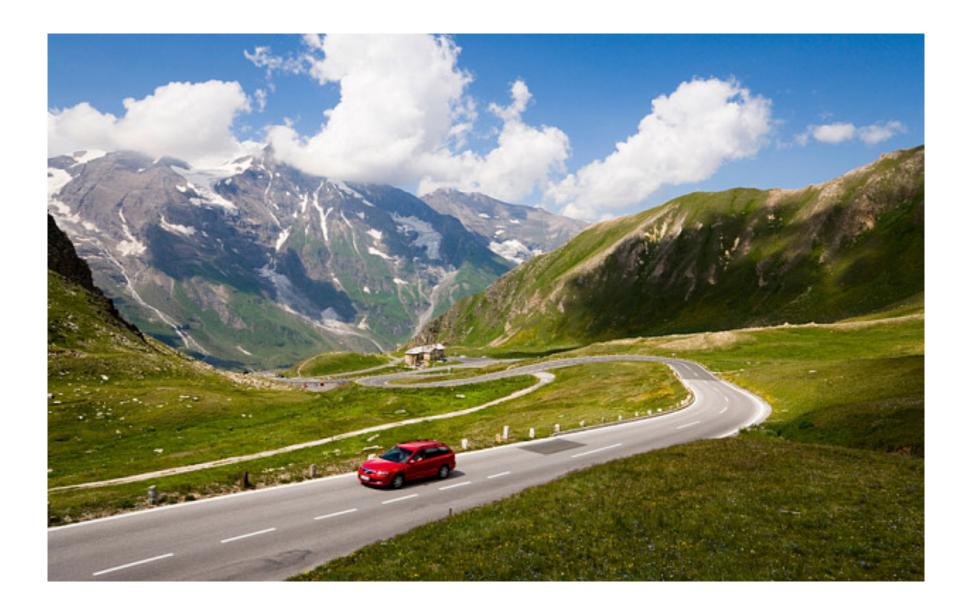


Hardware availability can change dramatically due to failure





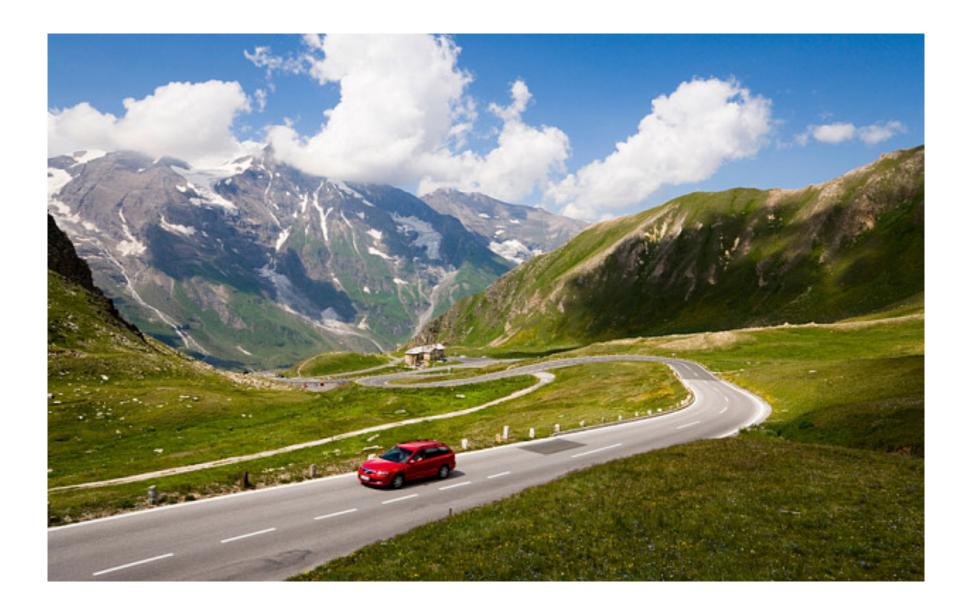
Hardware availability can change dramatically due to failure







Hardware availability can change dramatically due to failure









- Hardware availability can change dramatically due to failure
 - This is "normal" fault tolerance: how can software react to unexpected changes in the operating environment such as the failure of hardware
 - A wide range of techniques (too many to list)
 - But these are not the only adaptation scenarios to consider!





• Key hardware parameters vary across systems





Key hardware parameters vary across systems







• Key hardware parameters vary across systems









- Key hardware parameters vary across systems
 - Same architecture, but different specific properties (more cores, less memory, etc.)
 - Often can be handled by changing parameters in software (number of threads, parameters that govern cache usage, etc.)
 - Basic algorithms/implementations usually do not change
 - Often done at compile time or install time (hardware parameters like this do not usually change at run time)
 - Not always possible some scenarios require different parameters for different inputs, so you need to change at runtime
 - **One challenge:** how do you determine what to set a parameter to? lacksquare





- Key hardware parameters vary across systems
 - Same architecture, but different specific properties (more cores, less memory, etc.)
 - Loop tiling is a classic compiler technique to restructure a loop • Often (so its working set fits in cache. The tile size must be tuned for param different cache sizes
 - Bas Point blocking [Jo and Kulkarni 2011, 2012; Weijiang et al. 2015]
 - Often (is the equivalent of tiling but for recursive tree programs. Tuning usually happens at runtime because "tile size" is input dependent
 - Not always possible some scenarios require different parameters for different inputs, so you need to change at runtime
 - **One challenge:** how do you determine what to set a parameter to? lacksquare



f threads,

his do not































- Hardware platforms often change dramatically
 - Software that is written for one set of assumptions that hold true for one type of hardware may not hold true for a different type of hardware
 - Modern GPUs allow programmers to easily write parallel programs
 - But simply porting over programs that worked well on CPUs does not work GPUs have different execution models that must be managed to get good performance
 - Most work is compile-time or even earlier: compilers restructure code to work well on GPUs, or programmers use different design principles to write GPU code
 - One challenge: how can programmers make as few changes to their code as \bullet possible to have high performance GPUs





Hardware platforms often change dramatically

Our work looks at algorithmic and compiler techniques to map behavior) to GPUs effectively [Jo and Kulkarni, 2013; Goldfarb and Kulkarni 2013; Ren et al 2015; Liu et al 2016; Ren et al 2017] hot work t good We have also looked at targeting distributed memory machines in to work

- - GPL pert

• Software that is written for one set of assumptions that hold true for one type of hardware may not hold true for a different type of hardware Moder irregular programs (that have unpredictable control and data) • But Most v the same way [Hegde et al 2017a, b] well on GPUs, or programmers use different design principles to write GPU code

possible to have high performance GPUs



One challenge: how can programmers make as few changes to their code as

































elasticity

- Hardware resources change at runtime in a planned way
 - Hardware resources get added or removed at runtime in response to, e.g., changing demand (think: cloud computing)
 - Key difference from fault tolerance: hardware change is planned
 - Software must be designed to dynamically adapt to different hardware resources
 - Should take advantage of those resources (performance should be proportional to resources)
 - Should happen without interruption/disruption of software behavior (i.e., without restarting software)
 - One challenge: how can we write software so programmers do not have to consider the challenges of implementing elasticity





elasticity

Hardware resources change at runtime in a planned way

- - Sho transparently provides elasticity. prop

• Hardware resources get added or removed at runtime in response to, e.g., changing domand (thinks aloud computing) InContext [Yoo et al. 2011] and EventWave [Chuang et al. 2013] • Key provide programming models for *transparent elasticity*: Softwa programmers write *inelastic* software that does not consider dynamically changing hardware, and a runtime system resources Key challenge: how to restrict the programming model as little as (i.e.,

- Sho possible while providing transparent elasticity without restarting sortware
- consider the challenges of implementing elasticity



• One challenge: how can we write software so programmers do not have to







Unexpected (but normal) inputs may lead to unexpected outcomes



Unexpected (but normal) inputs may lead to unexpected outcomes









Unexpected (but normal) inputs may lead to unexpected outcomes











- Unexpected (but normal) inputs may lead to unexpected outcomes
 - Handling "out of the ordinary" inputs is part of the security problem, but what about inputs that a developer did not expect?
 - e.g., inputs much larger than expected?
 - If a developer has not tested against that kind of input, program may not do what they expect
 - Does a program still perform adequately? Does the program even run ۲ correctly at all?
 - **One challenge:** how do you know whether a program is behaving "incorrectly" when presented with an unexpected input?



input scaling



- - what about inputs that a developer did not expect?
 - e.g Vrisha [Zhou et al, 2011], Abhranta [Zhou et al, 2013], and
 - not do these models to determine if a program is behaving abnormally at large scales
 - \bullet correctly at all?
 - "incorrectly" when presented with an unexpected input?



input scaling

Unexpected (but normal) inputs may lead to unexpected outcomes

Handling "out of the ordinary" inputs is part of the security problem, but

 If a de of program behavior as a function of input scale, and then use m may

Does a program suit perform adequatery cloes the program even run

One challenge: how do you know whether a program is behaving



many dimensions to adaptability

- **Design principle or (semi) automatic technique?**
- When should software adapt?
 - Compile time vs. start up time vs. continuous
- How do we introduce adaptability?
 - Manual vs. compiler-driven vs. run-time



