

Software engineering for scientists

Paolo Ciancarini – paolo.ciancarini@unibo.it
Department of Informatics – University of Bologna

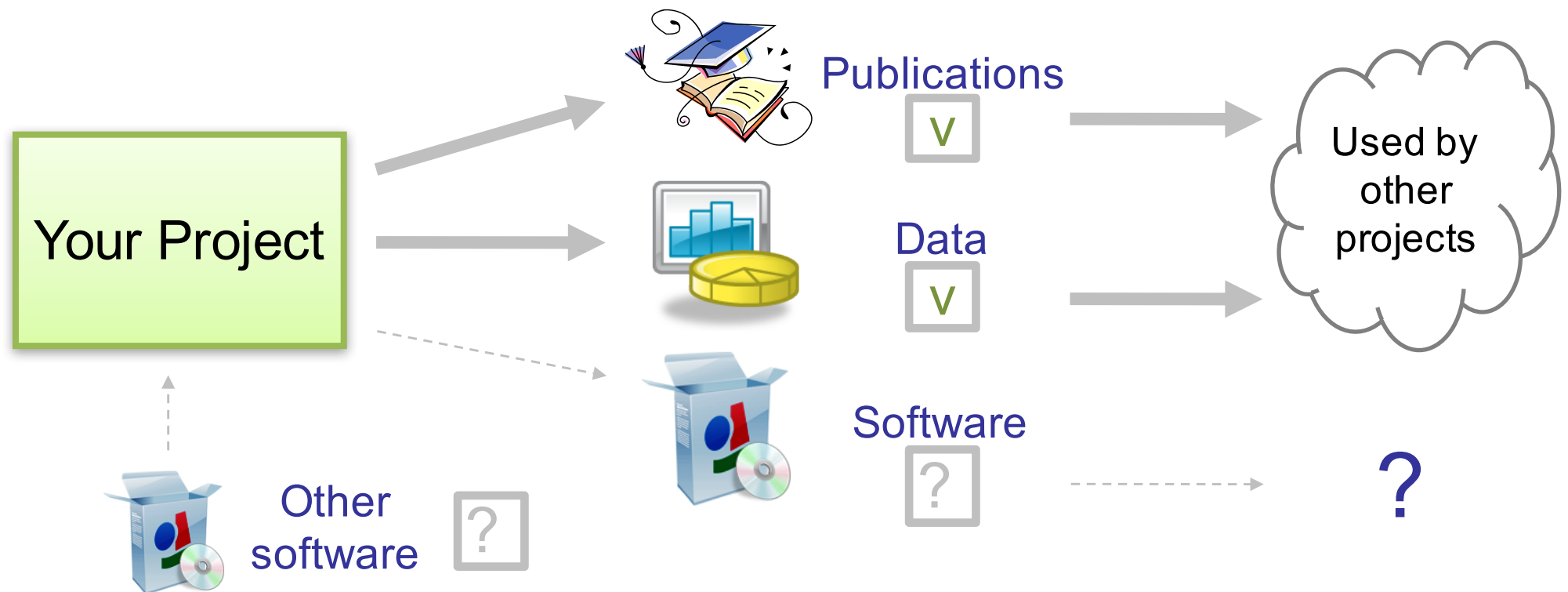
25th Summer School
on Parallel Computing

Bologna June 6, 2016

Agenda

- What is Software Engineering?
 - and why it is useful for scientists
- The software development lifecycle
- Sw development activities and tools
- Best practices for sw development

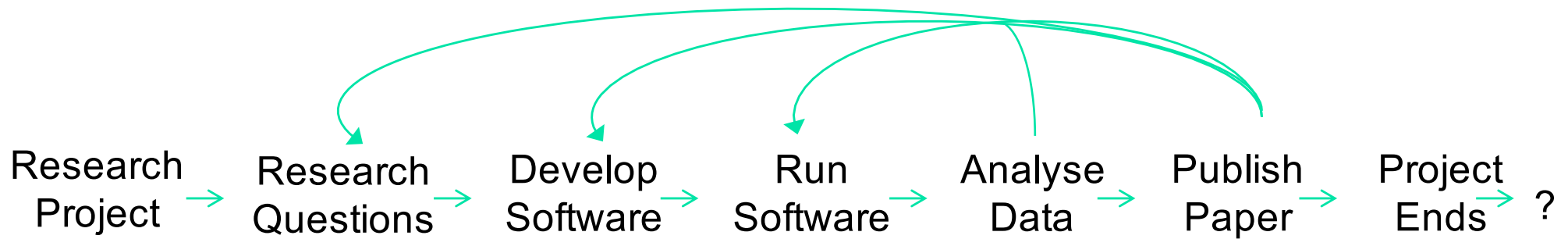
Using software in a research project



Do you agree?

- Scientists want to repeatedly tweak queries and analyses on their data sets and get immediate feedback. As long as you can clearly explain what you did to get the results published in a scientific paper, you don't need to pursue the code any further
- As science becomes increasingly computational in nature, it will become more important that scientific code does not end its **development cycle** on publication of the paper.
- Data exploration drives primary scientific discovery, but in order for future scientists to fully leverage the work of their predecessors, **robust**, **reproducible**, and **sustainable** software is needed to automate the parts we already know how to do

Typical development of software for science



What happens to the software?

- Thrown away
- Kept on some systems, possibly in different versions
- Dumped on a code repository

What happens when...

- You have a follow-on project?
- Someone wants to (re)use the code?
- Someone wants to reproduce your results?
- Maintenance or future reuse should be considered?

Software qualities

- **Robust** software: able to cope with errors during execution
- **Reproducible** software: version control
- **Sustainable** software: long lasting software able to cope with changes

Reproducibility?

- Tracking and recreating every step of your work
- Git: an enabling tool – use version control for everything
 - Paper writing
 - Grant writing
 - Everyday research
- Advantages:
 - A time machine view tracking every result
 - Distributed backup
 - Collaborate with colleagues

Sustainability?

For instance, software can *age*

- Ill-conceived design or modifications
- Functional operation degrades over time
- Lack of proper maintenance
- Infrastructure (os, libraries, language platform) evolves
- Some software types more susceptible

→ software becomes unsustainable, that is unusable

Software development: typical problems

- Unacceptable software performance
- Software hard to maintain or extend
- Inaccurate understanding of user needs
- Inability to deal with changing requirements
- Late discovery of serious flaws

→ poor software quality

What is software quality?

- Software *functional quality* reflects how well it complies with or conforms to a given design, based on some functional requirements
- Software *structural quality* refers to how it meets non-functional requirements that support the delivery of the functional requirements

IMPORTANT: Software quality can be measured!

Traditional approach

Sw development is a sequence including the following phases:

1. Requirements Analysis
2. Design
3. Coding
4. Testing: first check the units, then the system

The traditional development process goes through these phases **linearly**: first all the requirements are defined, then the design is completed, and finally the code is written and tested.

The key assumptions are that when design begins, requirements no longer change. When coding starts, the design ceases to change. etc.

NB: This “traditional” approach is sometimes called “waterfall development”

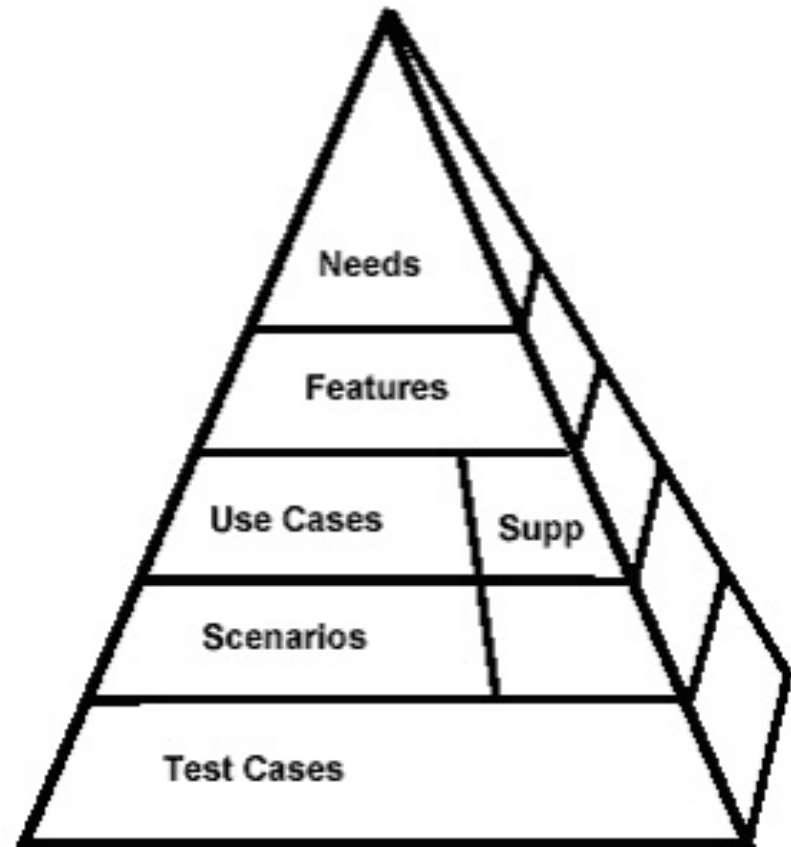
The requirements pyramid

Some user has some need

Needs are answered by “features” that some system must have

Each feature corresponds to a need and is a collection of requirements

Features and requirements can be aggregated in “scenarios”: after the code is built, testing it in the scenario will prove that its features satisfy the user’s needs



Poor practices for sw development

- Under-evaluation of development risks
- Overwhelming complexity
- Ambiguous communication
- Insufficient testing
- Insufficient requirements management
- Inconsistencies among requirements, designs, implementations, and tests
- Fragile software architecture

Risks when developing HPC sw

- Risks in sw development cycle
- Risks in the development environment
- Risks in the programming environment

See: Kendall, A Proposed Taxonomy for Software Development Risks for High-Performance Computing (HPC) Scientific/Engineering Applications, SEI 2007

Examples

- Typical risk in the development cycle: misunderstanding requirements
- Typical risk in the development environment: too many manual activities
- Typical risk in the programming environment: underestimating dependencies

Best practices for scientific computing

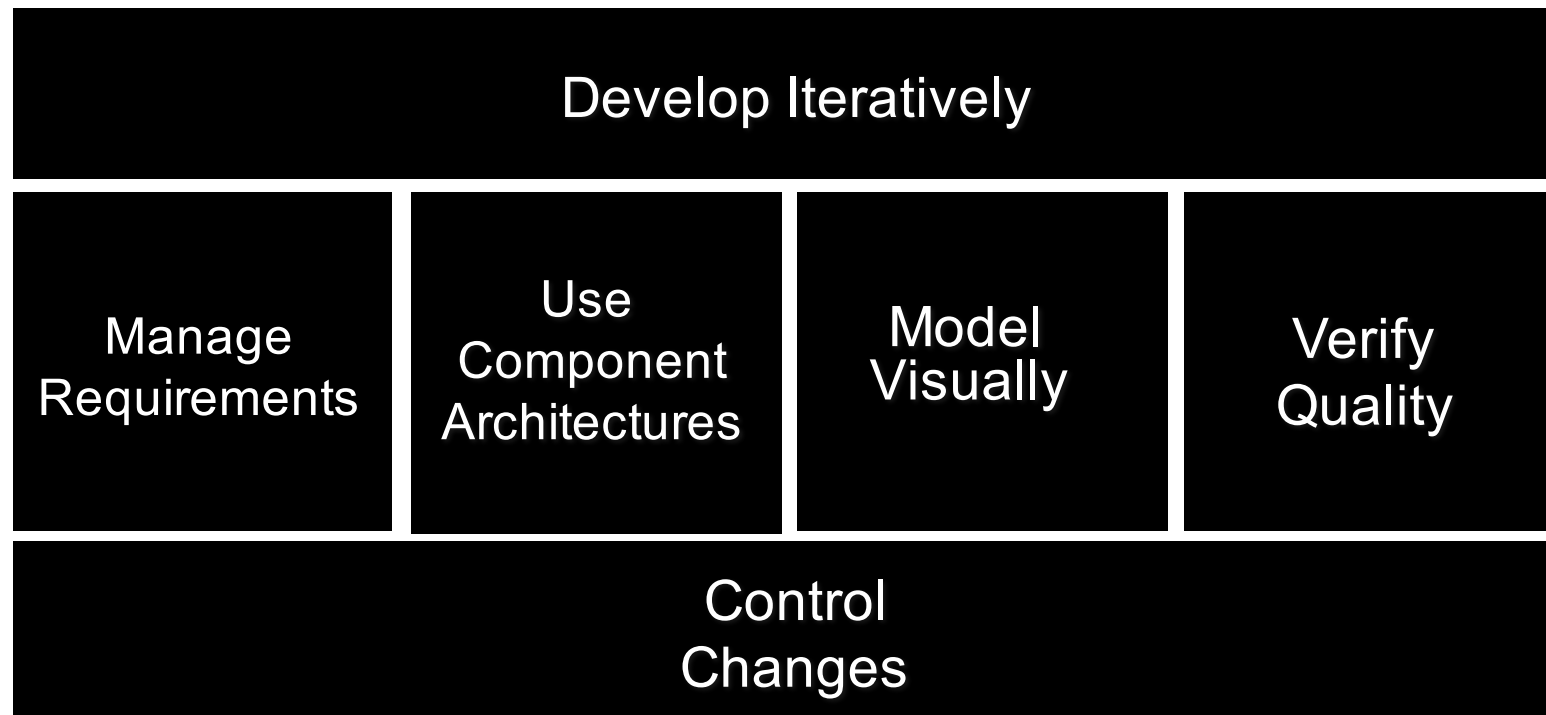
- Write programs for people, not computers
- Use a tool to automate workflows
- Make incremental changes, use a version control system
- Reuse code instead of rewriting it
- Plan for finding mistakes
- Optimize software only after it works correctly
- Document design and purpose (not mechanics)
- Collaborate (eg. by pair programming or by using an issue tracking tool)

Best practices of software development

- Develop iteratively
- Control changes

- Manage requirements
- Verify quality
- Use components
- Model software architecture visually

Best practices of software development



Know these!

Enters Software Engineering

“Software engineering is the discipline concerned with all aspects of software production from the early stages of system specification to maintaining the system after it has gone into use”

[Sommerville 2007]

Software Engineering

“The establishment and use of sound engineering **principles** in order to obtain **economically** software that is **reliable** and works efficiently on real machines.” [Naur & Randell, 1968]

Software Engineering

- A definition and some issues
 - “developing quality software on time and within budget”
- Trade-off between a system perfectly engineered and the available resources
 - SwEng has to deal with real-world issues
- State of the art
 - Community decides on “best practices” + life-long education

What is Software Engineering?

A naive view:



But ...

- Where did the *problem specification* come from?
- How do you know the problem specification corresponds to and satisfies the *user's needs*?
- How did you decide how to *structure* your program?
- How do you know the program actually *meets the specification*?
- How do you know your program will always *work correctly*?
- What do you do if the users' *needs change*?
- How do you *divide tasks up* if you have more than a one person in the developing team?
- How do you *reuse* existing software for solving similar problems?

What is Software Engineering?

“multi-person construction of multi-version software”

— Parnas

- Software is complex and difficult to build
- Team-work
 - Scale issue (“program well” is not enough) + communication issues: Conway’s law
- Successful software systems must evolve or perish
 - Changes to the software is the norm, not the exception

Conway's Law

- The law: *Organizations that design software systems are constrained to produce designs that are copies of the communication structures of these organizations*
- Example: "If you have four groups working on a compiler, you'll get a 4-pass compiler"
- Several studies found significant differences in modularity when software is outsourced, consistent with a view that distributed teams tend to develop more modular products

What is Software Engineering?

“software engineering is different from other engineering disciplines”

— Sommerville

- It is not constrained by physical laws
 - limit = human knowledge
- It is constrained by social forces
 - Balancing stakeholders needs
 - Consensus on functional and especially non-functional requirements

Requirements

Software design

Coding

Development process

Testing

Software Engineering

Evolution

Sw quality

IPR & licensing

Tools

Project management

Configuration management

Topics of the discipline

- Product and process standards for software
- Project management for software systems
- Software development models: planned vs agile
- Requirement analysis
- Software design by visual modeling
- Measuring, verifying, and ensuring software quality
- Software evolution and maintenance
- Typical tools used by software engineers:
 - Version control, configuration management

Software Engineering for HPC

- Software engineering aims to designing, implementing, and modifying software so that it is faster to build, of higher quality, more maintainable
- In HPC we have all the common problems of software development, plus the specific problem that software developers have scarce knowledge of software engineering best practices
- In the following slides we will deal with some of these problems and suggest some solutions

Roadmap

- What is Software Engineering?
- **The Software Development Lifecycle**
- Software development activities
- Methods and tools

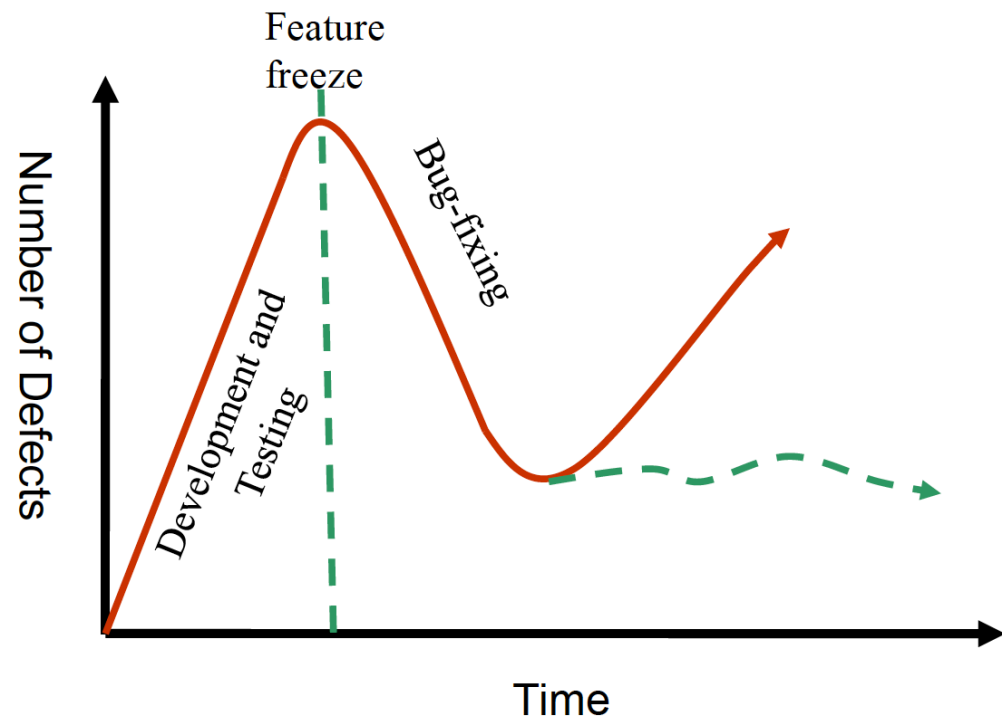
Software: the product of a process

- Many kinds of software products
 - many kinds of development processes
- “Study the process to improve the product”
- A software development process can be described according to some specific “model”
- Examples of process models: waterfall, iterative, agile, explorative,...
- These models differ mainly in the roles and activities that the stakeholders cover

Beware of software aging!

Software can *age*

- Ill-conceived design or modifications
- Functional operation degrades over time
- It becomes unsustainable, unusable
- Lack of proper maintenance
- Infrastructure (os, libraries, language platform) evolves
- Some software types more susceptible



Stakeholders

Typical stakeholders in a sw process

- Users
- Decisors
- Designers
- Management
- Technicians
- Funding people
- ...

Each stakeholder has a specific viewpoint on the product and its development process

Just a joke?



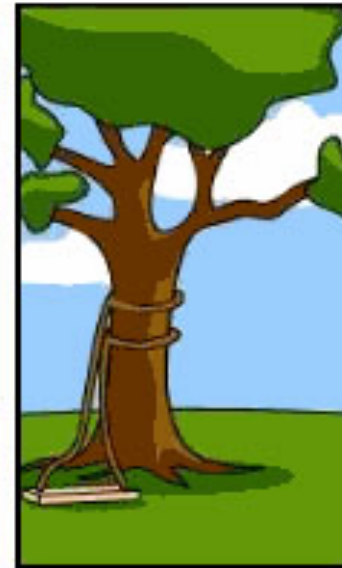
How the customer explained it



How the Project Leader understood it



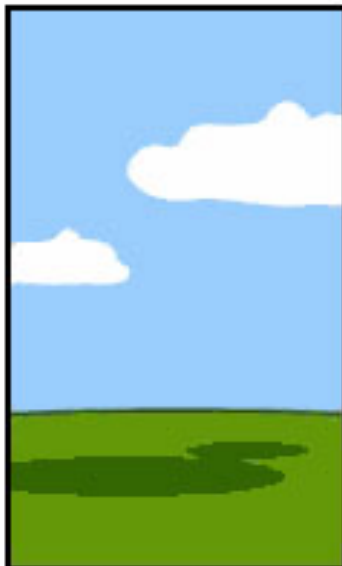
How the Analyst designed it



How the Programmer wrote it



How the Business Consultant described it



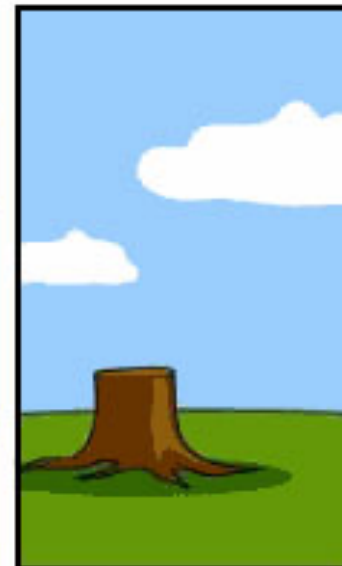
How the project was documented



What operations installed



How the customer was billed



How it was supported



What the customer really needed

HPC stakeholders attributes

Attribute	Values	Description
Team size	Individual	This scenario, sometimes called the “lone researcher” scenario, involves only one developer.
	Large	This scenario involves “community codes” with multiple groups, possibly geographically distributed.
Code life	Short	A code that’s executed few times (for example, one from the intelligence community) might trade less development time (less time spent on performance and portability) for more execution time.
	Long	A code that’s executed many times (for example, a physics simulation) will likely spend more time in development (to increase portability and performance) and amortize that time over many executions.
Users	Internal	Only developers use the code.
	External	The code is used by other groups in the organization (for example, at US government labs) or sold commercially (for example, Gaussian, www.gaussian.com)
	Both	“Community codes” are used both internally and externally. Version control is more complex in this case because both a development and a release version must be maintained.

V.Basili et al., Understanding the High-Performance-Computing Community: A Software Engineer’s Perspective, IEEE Software, 2008

The software development process

- **Software process**: set of roles, activities, and artifacts necessary to create a software product
- Example **roles**: stakeholder, designer, developer, tester, maintainer, ecc.
- Example **artifacts**: source code, libraries, comments, test suites, etc.

Activities

- Each organization differs in the products it builds and the way it develops them; however, most development processes include:
 - Specification
 - Design
 - Verification and validation
 - Evolution
- The development activities must be modeled to be managed and supported by automatic tools

Software development activities

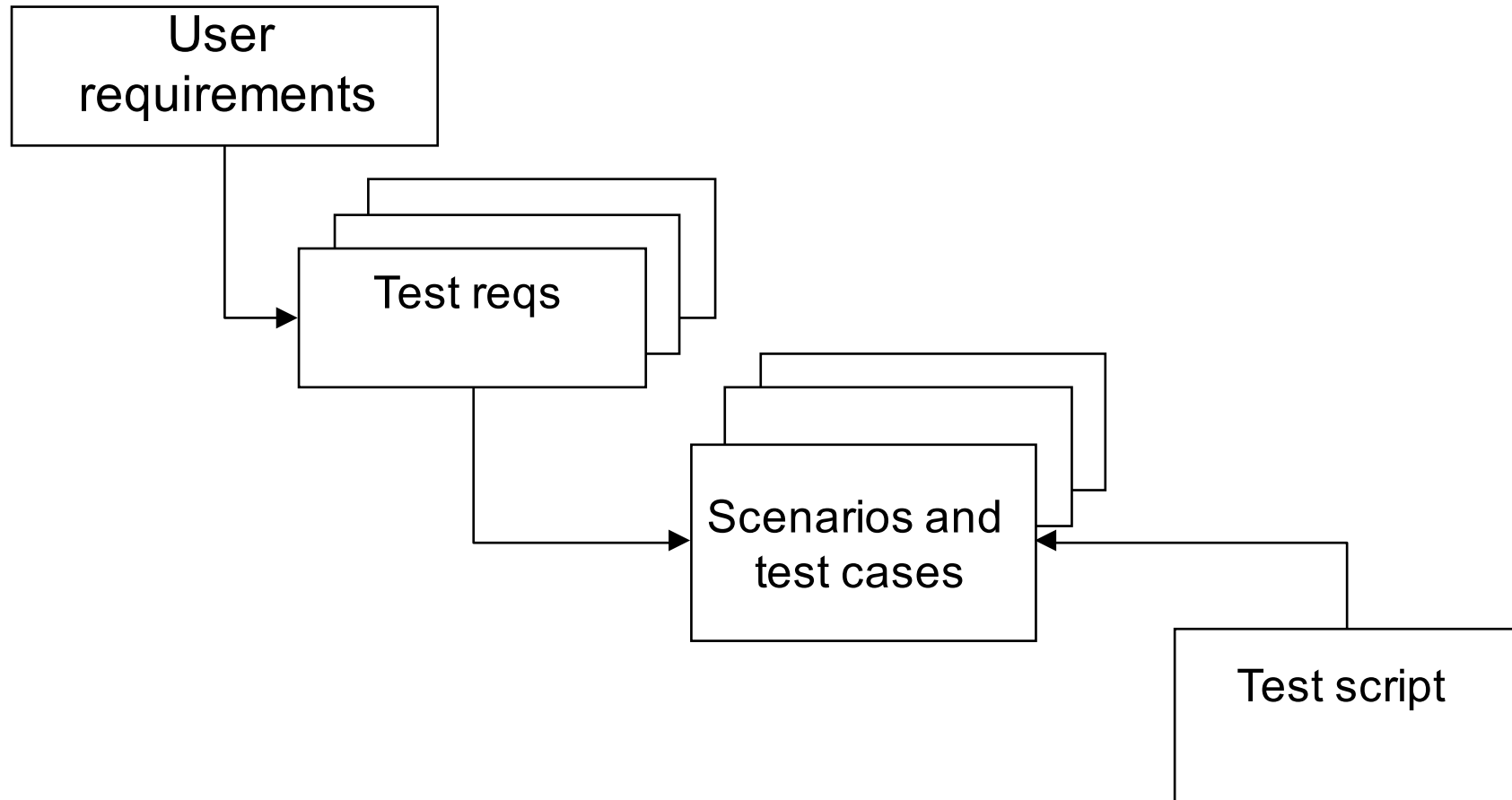
<i>Requirements Collection</i>	Establish customer's needs
<i>Analysis</i>	Model and specify the requirements ("what")
<i>Design</i>	Model and specify a solution ("how")
<i>Implementation</i>	Construct a solution in software
<i>Testing</i>	Validate the software against its requirements
<i>Deployment</i>	Making a software available for use
<i>Maintenance</i>	Repair defects and adapt the sw to new requirements

NB: these are ongoing activities, not sequential phases!

First development step: requirements

- The first step in any development process consists in understanding the needs of someone asking for a software
- The needs should be stated explicitly in “requirements”, which are statements requiring some function or property to the final software system

Requirements and tests



Types of testing

Acceptance testing (by the user)

Performance testing

System testing

Integration testing

Unit testing

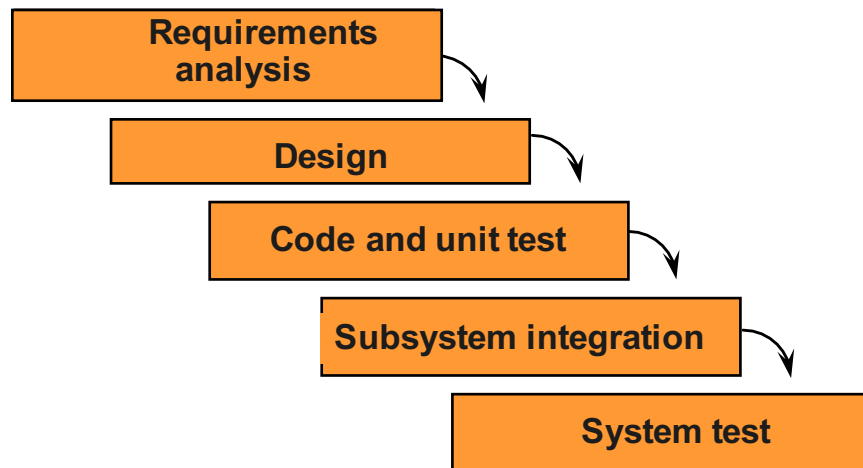
Models for the software process

A model for the software development process is a method to describe the roles, the tasks, and the documents to be developed

- Waterfall (planned, linear)
- Spiral (planned, iterative)
- Agile (unplanned, test driven)

Waterfall characteristics

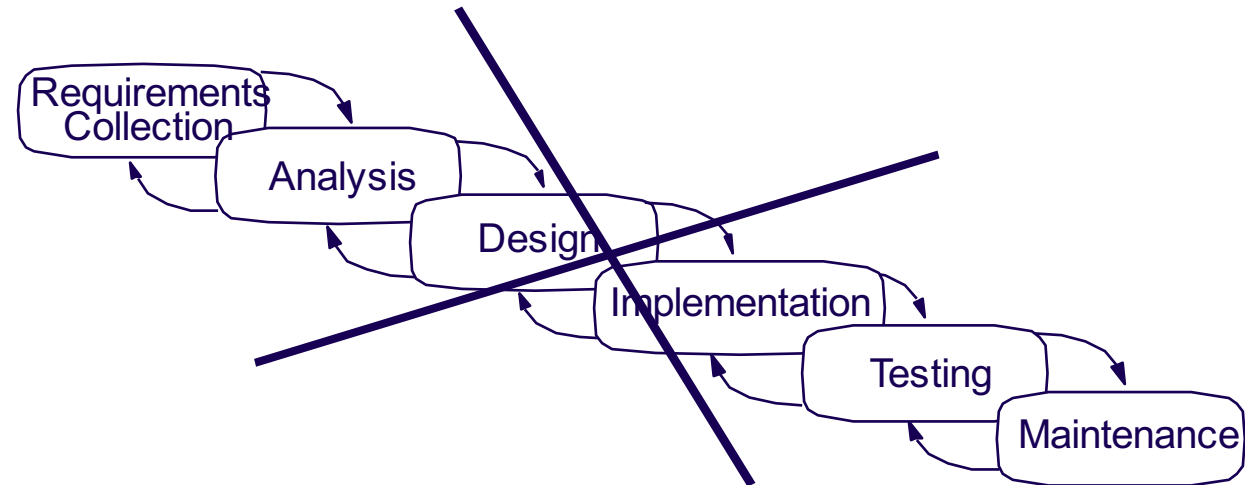
Waterfall Process



- One way communications
- Delays confirmation of critical risk resolution
- Measures progress by assessing work-products that are poor predictors of time-to-completion
- Delays and aggregates integration and testing
- Precludes early deployment
- Frequently results in major unplanned iterations

The classical software lifecycle

The classical software lifecycle models the software development as a step-by-step “waterfall” between the various development phases



The waterfall model is flawed for many reasons:

- Requirements must be ***frozen too early*** in the life-cycle
 - User requirements are ***validated too late***
- **Risks** in constructing wrongly the software are high

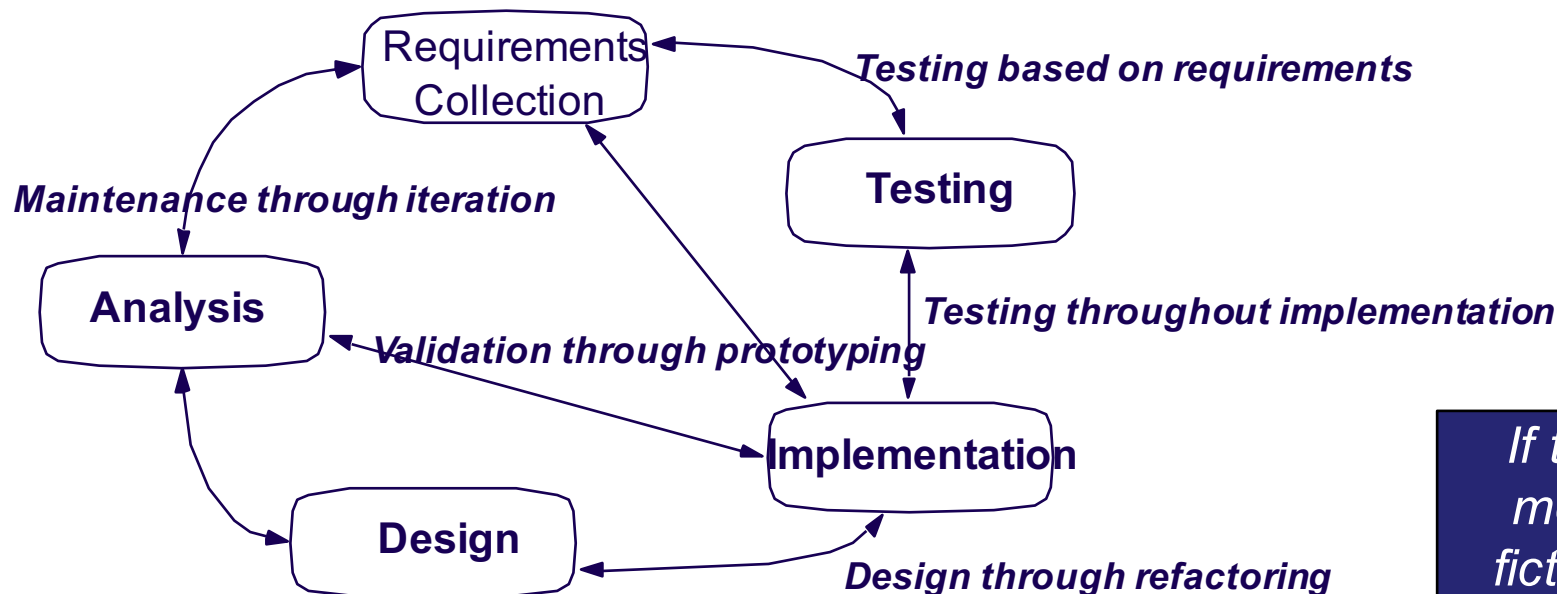
Problems with the waterfall lifecycle

1. “Real projects rarely follow the sequential flow that the waterfall model proposes. *Iteration* always occurs and creates problems in the application of the paradigm”
2. “It is often *difficult* for the customer *to state all requirements* explicitly. The classic life cycle requires this and has difficulty accommodating the natural uncertainty that exists at the beginning of many projects.”
3. “The customer must have patience. A *working version* of the program(s) will not be available until *late in the project* timespan. A major blunder, if undetected until the working program is reviewed, can be disastrous.”

— Pressman, SE, p. 26

Iterative development

In practice, development is always iterative,
and *most* activities can progress in parallel



If the waterfall model is pure fiction, why is it still the dominant software process?

Iterative development

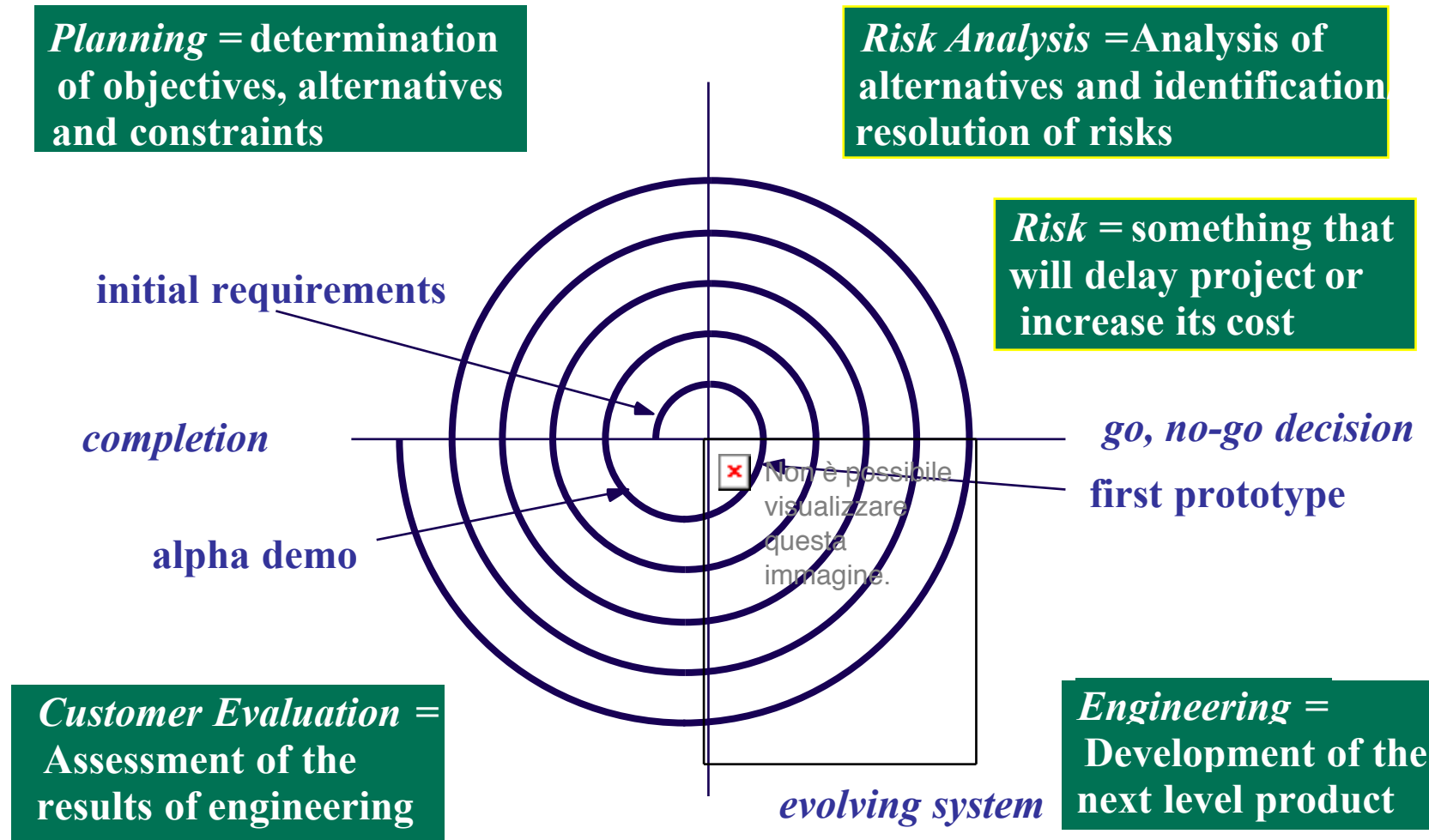
- Plan to iterate your analysis, design and implementation
 - You will not get it right the first time, so integrate, validate and test as frequently as possible
 - During software development, more than one iteration of the software development cycle may be in progress at the same time
 - This process may be described as an 'evolutionary acquisition' or 'incremental build' approach

Iterative development

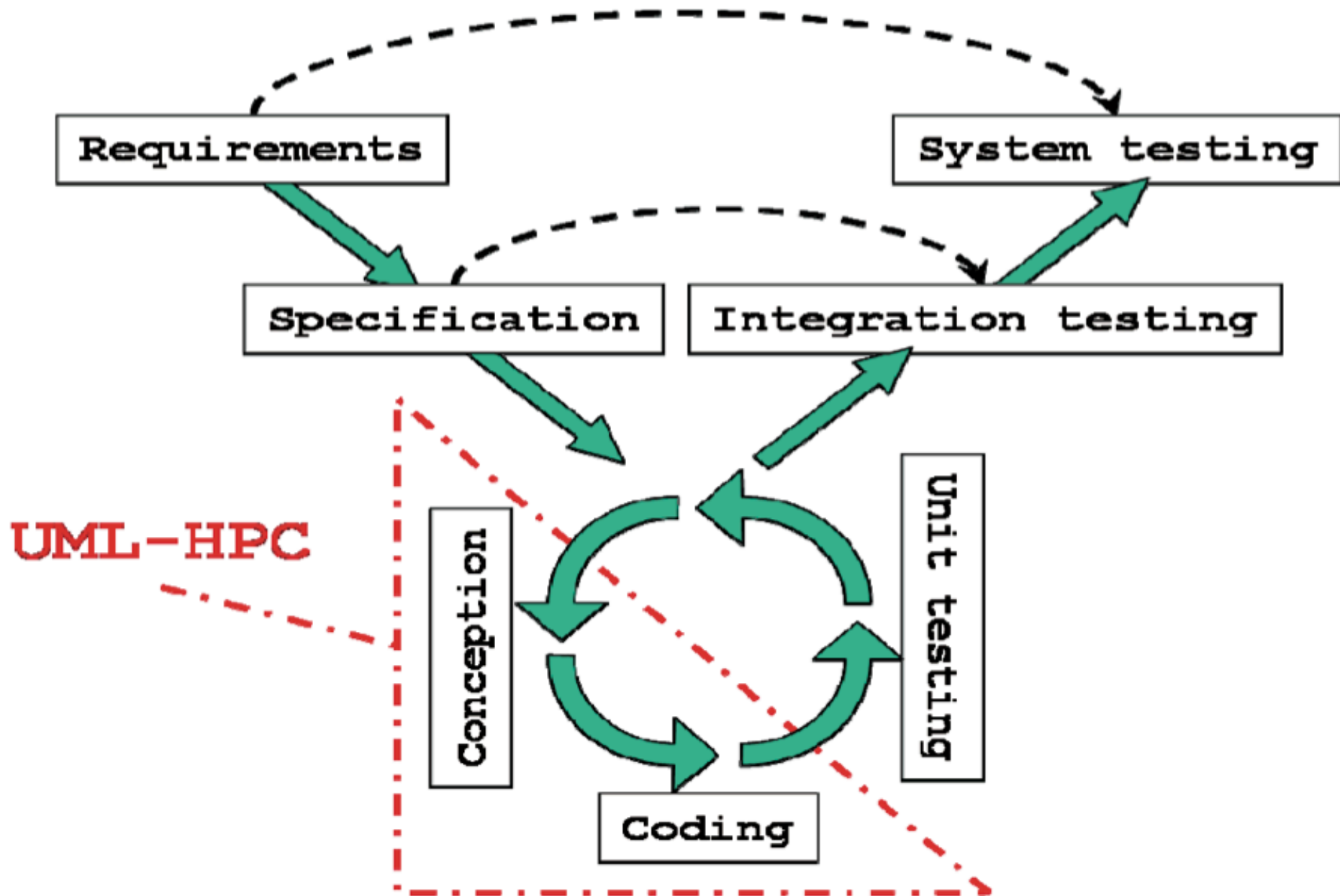
Plan to *incrementally* develop (i.e., prototype) the system

- If possible, *always have a running version* of the system, even if most functionality is yet to be implemented
- *Integrate* new functionality as soon as possible
- *Validate* incremental versions against user requirements.

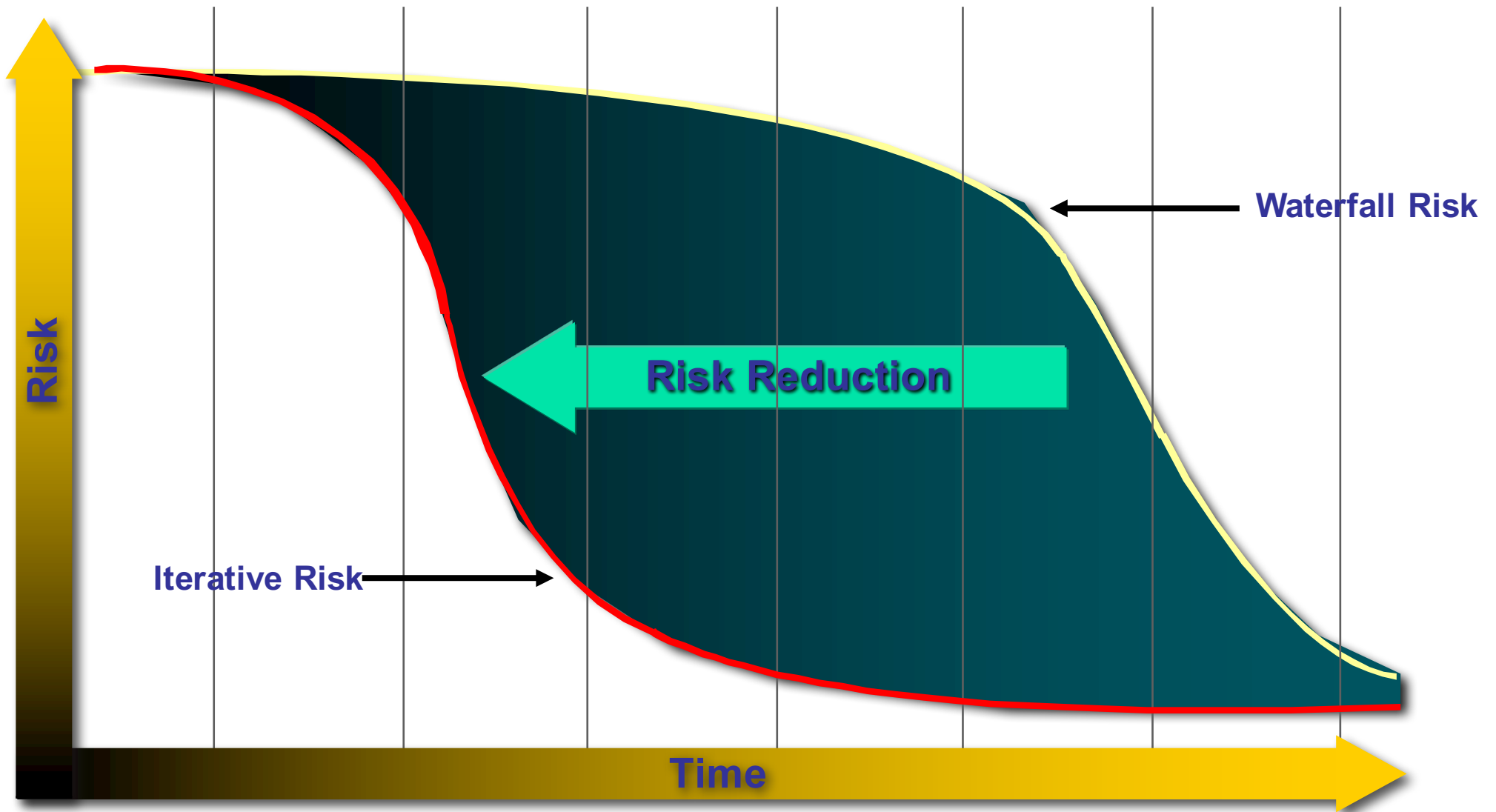
The spiral lifecycle



A process for HPC [Lugato 2010]

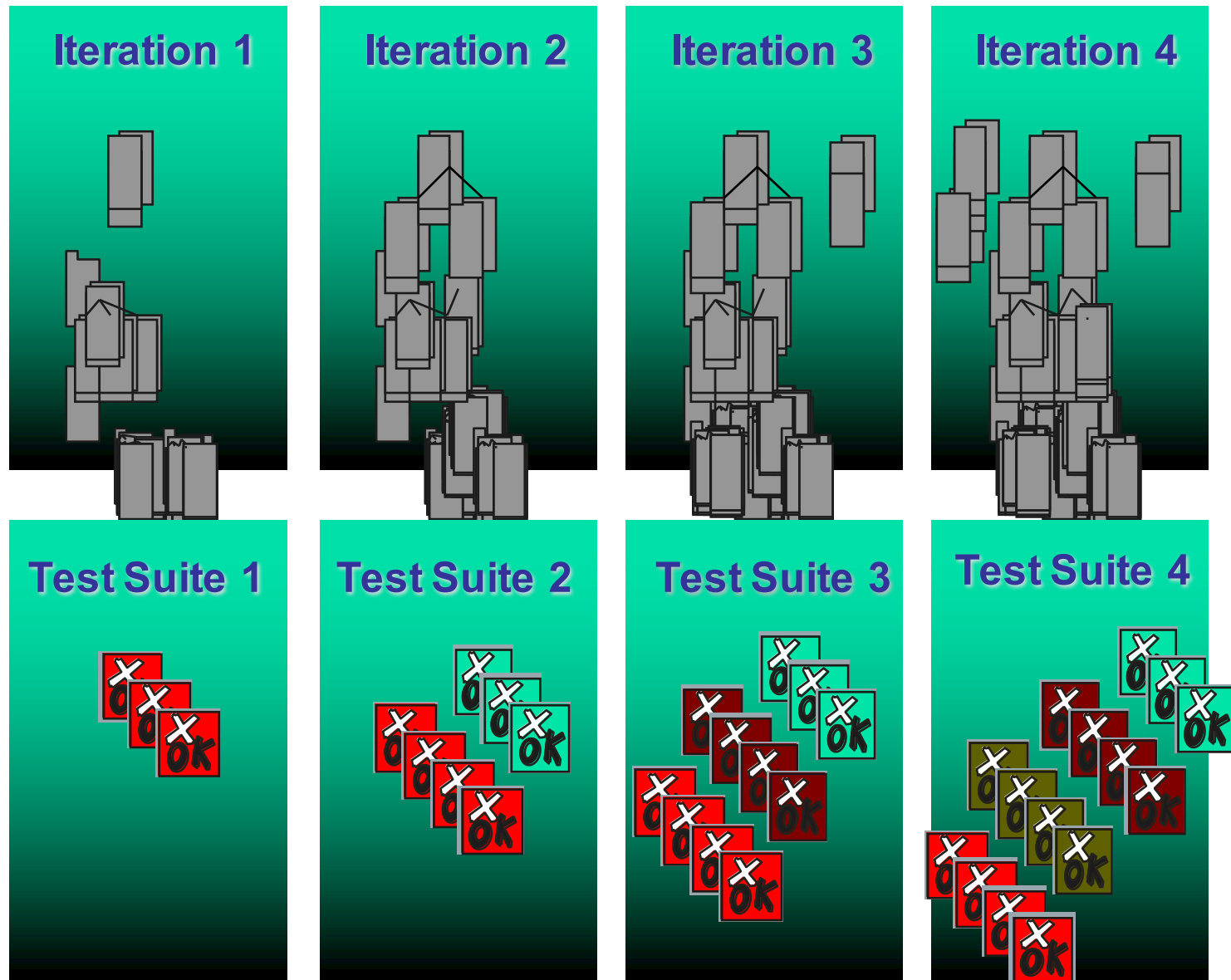


Risk: waterfall vs iterative



Test each iteration

Requirements,
models
and code



Tests

Testing before designing

- What is software testing? an investigation conducted to provide information about the quality of some software product
- In planned process models testing happens after the coding, and checks if the code satisfies the requirements
- What happens if we define the tests **before** the code they have to investigate?

Agile development processes

- There are many agile development methods; most minimize risk by developing software in short amounts of time
- The requirements are initially grouped in stories and scenarios
- Then the tests for each scenario are agreed with the user, before any code is written
- Each code is tested against its scenario tests, and integrated after it passes its unit tests

Agile ethics

- www.agilemanifesto.org

We are uncovering better ways of developing software by doing it and helping others do it.
Through this work we have come to value:

Individuals and interactions over **processes and tools**
Working software over **comprehensive documentation**
Customer collaboration over **contract negotiation**
Responding to change over **following a plan**

That is, while there is value in the items on **the right**, we prefer the items on **the left**.

- Management can tend to prefer the things on the right over the things on the left

SCRUM

Team-Level Planning

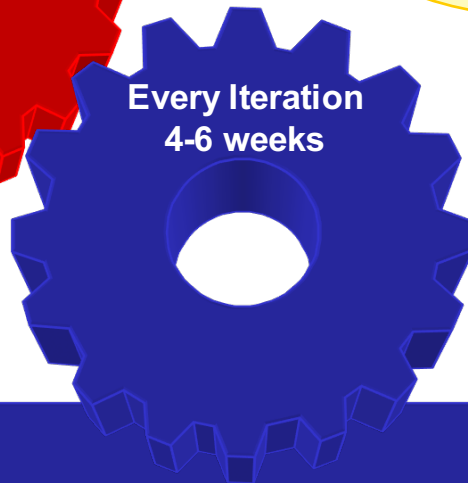


Daily Scrum Meeting:

15 minutes

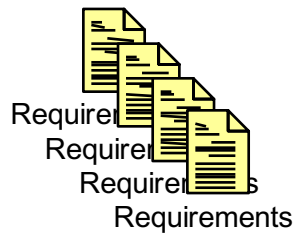
Each team member answers 3 questions:

- 1) What did I do since last meeting?
- 2) What obstacles are in my way?
- 3) What will I do before next meeting?

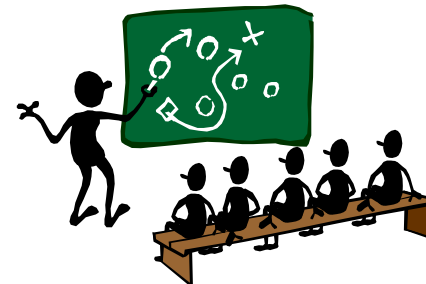


Working Software
Delivered

Prioritised
Iteration
Scope



Prioritised Requirements &
Features "Backlog"



Applying Agile:

Continuous integration; continuously monitored progress

Roadmap

- What is Software Engineering?
- The Software Development Lifecycle
- **Software Development Activities**
- Methods and tools

Requirements collection

User requirements are often expressed *informally*:

- They are grouped in *features*
- They are put in context in usage scenarios

Even if requirements are documented in written form, they may be *incomplete*, *ambiguous*, or *incorrect*

Changing requirements

Requirements *will* change!

- *inadequately captured* or expressed in the first place
- user and business *needs may change* during the project

Validation is needed *throughout* the software lifecycle, not only when the “final system” is delivered!

- build constant *feedback* into your project plan
- plan for *change*
- early *prototyping* [e.g., UI] can help clarify requirements

Requirements analysis

Analysis is the process of specifying *what* a system will do

- The goal is to provide an understanding of what the system is about and what its underlying concepts are

The result of analysis is a *specification document*

Does the requirements specification correspond to the users' actual needs?

Design

Design is the process of specifying *how* the specified system behaviour will be realized from software components. The results are *architecture* and *detailed design documents*.

Object-oriented design delivers models that describe:

- how system operations are implemented by *interacting objects*
- how classes refer to one another and how they are related by *inheritance*
- *attributes* and *operations* associated to classes

*Design is an iterative process,
proceeding in parallel with
implementation!*

Implementation and testing

Implementation is the activity of *constructing* a software solution to the customer's requirements.

Testing is the process of *verifying* that the solution meets the requirements.

- The result of implementation and testing is a *fully documented* and *verified* solution.

What Is Software Testing?

- Testing is the process of executing a program with the intent of finding errors
 - Assume that the program contains errors and then test the program to find as many of the errors as possible
- Testing is a destructive - creative process which aims at establishing confidence that a program does what it is supposed to do

Testing!

1

- Provide automated build process
 - Far easier & quicker to validate changes
 - e.g. Make, Ant, Maven

2

- Provide automated regression test suite - TDD
 - Do changes break anything?
 - JUnit, CPPUNIT, xUnit, fUnit, ...

3

- Join together: automated build & test
 - A 'fail-fast' environment

4

- Infrastructure support
 - Nightly builds – run build & test overnight, send reports
 - Continuous integration - run build & test when codebase changes

Towards *anytime releasable code!*

Iterativity of design, Implementation and testing

Design, implementation and testing are iterative activities

- The implementation does not “implement the design”, but rather the design document *documents the implementation!*
- System tests reflect the requirements specification
- Testing and implementation go hand-in-hand
 - Ideally, test case specification *precedes* design and implementation

Maintenance

Maintenance is the process of changing a system after it has been deployed.

- *Corrective maintenance*: identifying and repairing defects
- *Adaptive maintenance*: adapting the existing solution to new platforms
- *Perfective maintenance*: implementing new requirements
- *Preventive maintenance*: repairing a software product before it breaks

In a spiral lifecycle, everything after the delivery and deployment of the first prototype can be considered “maintenance”!

Maintenance activities

“Maintenance” entails:

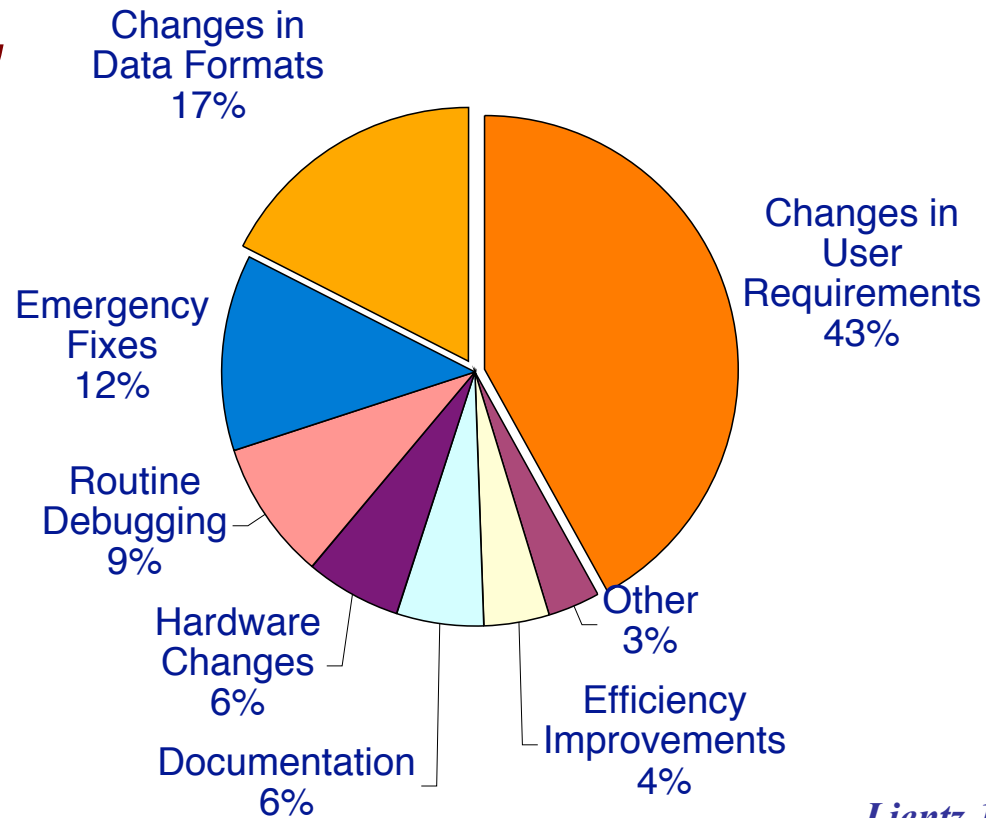
- configuration and version management
- reengineering (redesigning and refactoring)
- updating all analysis, design and user documentation

Repeatable, automated tests enable evolution and refactoring

Maintenance costs

“Maintenance”
typically accounts for
70% of software costs!

*Means: most
project costs
concern continued
development after
deployment*



– Lientz 1979

Deployment

- Virtual Machines
 - Software pre-installed, ready to run
 - Often easiest
 - Not enough in itself – documentation!
- Release software
 - Prioritise & select requirements -> Develop -> Test -> Commit changes to repository -> Test -> Release
 - Documentation (minimum: quick start guide)
- Licencing
 - Specify rights for using, modifying and redistributing

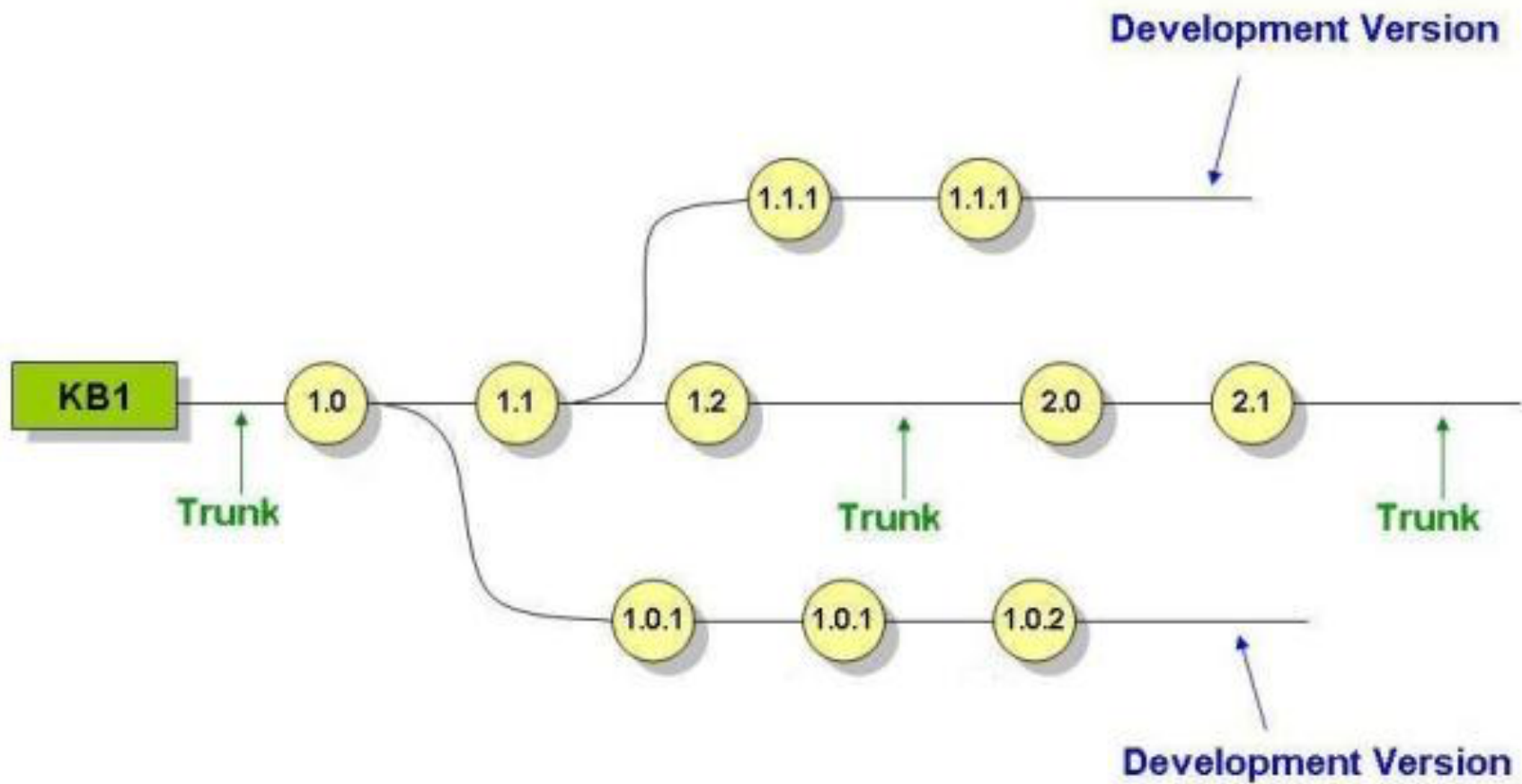
Configuration management

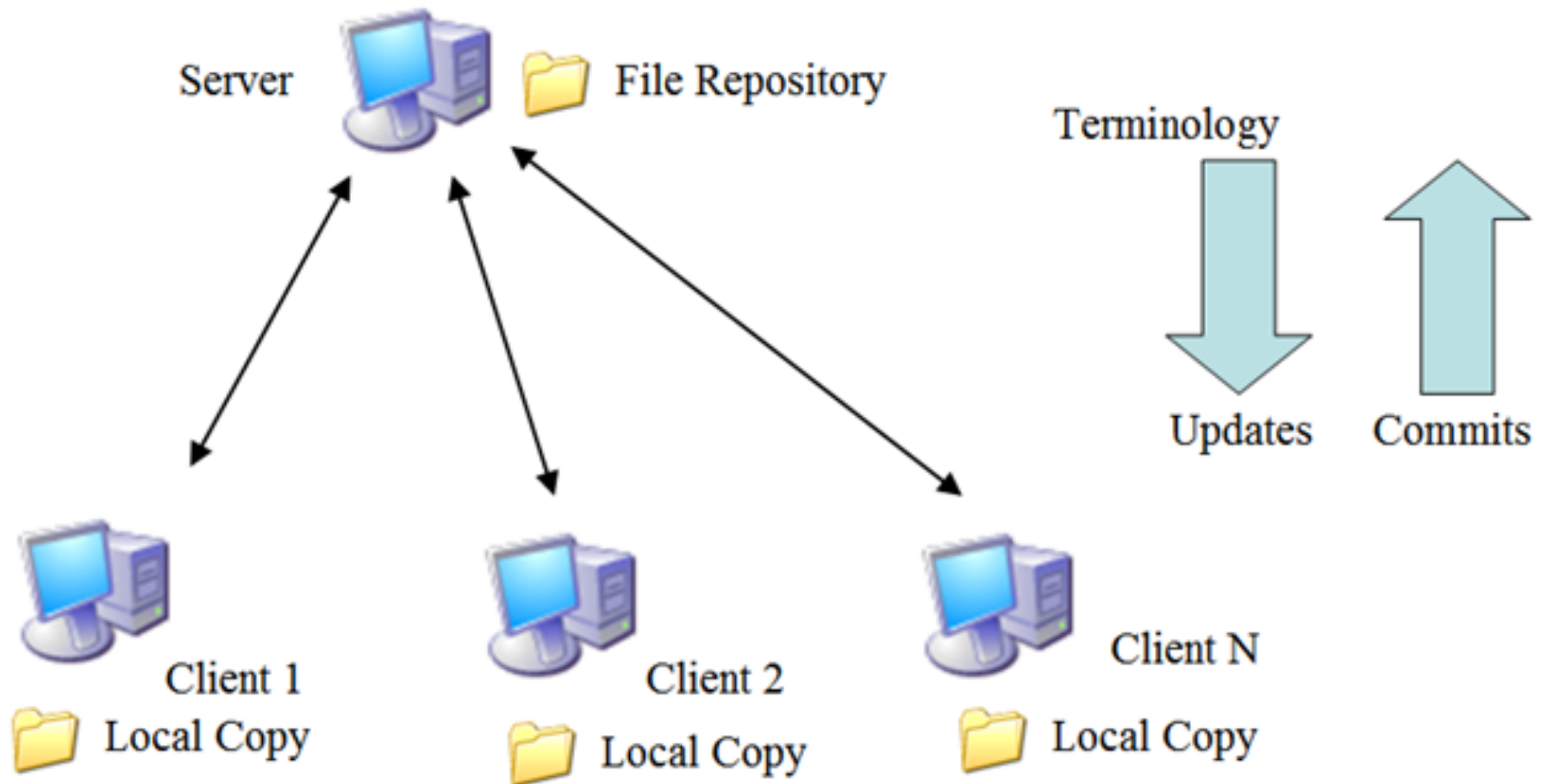
- Run your own CM system, if you have the resources
 - Generally easy to set up
 - Full control, but be sure to back it up!
- Some public solutions can offer most of these for free
 - SourceForge, GoogleCode, GitHub, Codeplex, Launchpad, Assembla, Savannah, ...
 - BitBucket for private code base (under 5 users)
 - See (for hosting code and related tools)
<http://software.ac.uk/resources/guides/choosing-repository-your-software-project>
 - See (for hosted continuous integration)
<http://www.software.ac.uk/blog/2012-08-09-hosted-continuous-integration-delivering-infrastructure>

“If you’re not using version control, whatever else you might be doing with a computer, it’s not science” – Greg Wilson, Software Carpentry

Version control

- Version management allows you to control and monitor changes to files
 - What changes were made?
 - Revert to previous versions
 - When were changes made ?
 - What code was present in release 2.7?
- Earliest tools were around 1972 (SCCS)
- Older tools – RCS, CVS, Microsoft Source Safe, PVCS Version Manager, etc...
- Current tools – Subversion, Mercurial, Git, Bazaar





Local

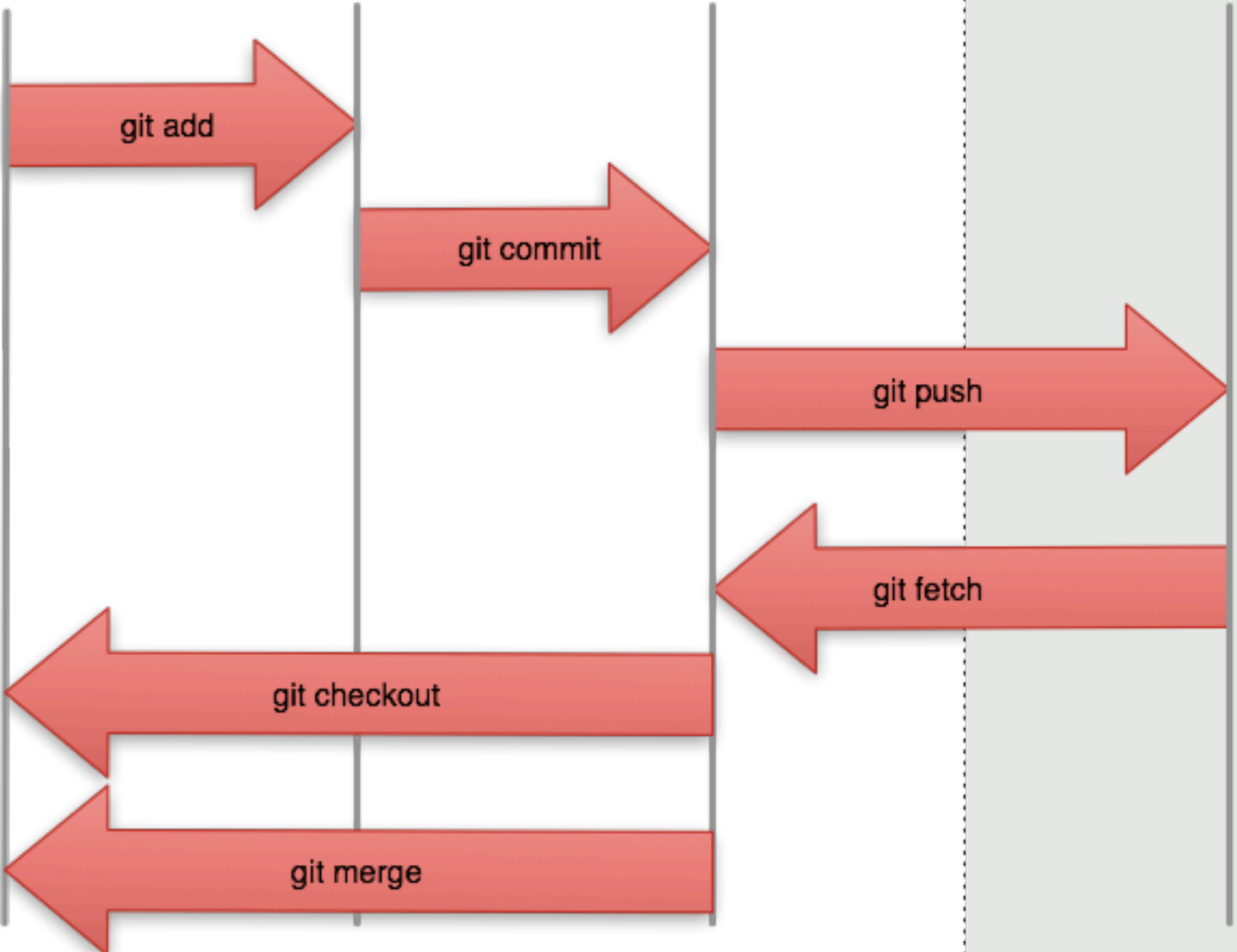
Remote

working directory

staging area

local repo

remote repo



Version control concepts

- checkout – get a local copy of the files
 - I have no files yet, how do I get them?
- add – add a new file into the repository
 - I created a new file and want to check it in
- commit – send locally modified files to the repository
 - I made changes, how do I send them to the group?
- update – update all files with latest changes
 - Other people made changes, how do I get them?
- tag / branch – label a “release”
 - I want to “turn in” a set of files

Conclusions

Software engineering deals with

- the way in which software is made (process),
- the languages to model and implement software,
- the tools that are used, and
- the quality of the result (testing and measures)

Self test questions

- How does Software Engineering differ from programming?
- Why is the “waterfall” model unrealistic?
- What is the difference between analysis and design?
- Why plan to iterate? Why develop incrementally?
- Why is programming only a small part of the cost of a “real” software project?

Reference: papers

- V.Basili et al., Understanding the High-Performance- Computing Community: A Software Engineer's Perspective, *IEEE Software*, 2008
- G. Wilson et al., Best Practices for Scientific Computing. *PLoS Biol* 12(1), 2014
- Kendall et al., A Proposed Taxonomy for Software Development Risks for High-Performance Computing (HPC) Scientific/Engineering Applications, TN-039 CMU, 2007
- D.Lugato et al., Model-driven engineering for HPC applications, *Proc. Modeling Simulation and Optimization Focus on Applications*, Acta Press (2010): 303-308.

References: books

- Pressman, *Software engineering a practitioner approach*, 8th ed., McGrawHill, 2014
- The Computer Society, *Guide to the Software Engineering Body of Knowledge*, 2013
www.computer.org/portal/web/swebok

Useful references

- software.ac.uk Software Sustainability Institute
- software.ac.uk/resources/case-studies
- software-carpentry.org Software carpentry
- Proc. 2015 Int. Workshop on Sw Engineering for HPC in Science

Questions?

