



Lecture 20: Build systems and dependencies

CS 5150, Spring 2026

Administrative Reminders

- **Final Presentation:** Apr 30/May 1. See the schedule on the shared sheet. Please be in the room at the start of each session.
- All members must be in person
- Client meeting scores of sprint 2 are shared, sprint 3 will be shared soon
- Peer Review 3 form will be shared soon, please submit on time

What should be in the final presentation?

- What did you plan to build? Context, goals
- What did you achieve so far? What was missed?
- How do you know if you built it right?
- Software engineering practices
- AI usage: what tools did you use? How useful were those?
(mandatory)
- User Testing
- Demo

Previously ...

Internal vs. External dependencies

Internal

- Maintainers' goals are (hopefully) aligned
- Can audit for all uses of a library
- Can coordinate large-scale changes of all code using library (facilitated by monorepo)
- Can manage with [source control](#) tools, policies

External

- Cannot assume coordination between library and users
- Cannot enforce compatibility, maintenance policies
- Cannot control release schedule
- Danger of **diamond dependency** problem
- Domain of [dependency management](#)

Dependency Management Practices

- **Version pinning:** select the exact dependency version
 - But only applies to direct dependencies
 - Problems?
- Signature and hash verification
- **Lockfiles:** pinning+sig/hash verification for full dependency tree
 - Compiles all dependencies and sub-dependencies (entire dependency tree)
 - Better reproducibility and consistency
- **Dependency confusion attack:** publishing projects with the same name as an internal project to open-source
- **Vulnerability scanning:** scan lock files to check the artifact versions

<https://cloud.google.com/blog/topics/developers-practitioners/best-practices-dependency-management>

Supply Chain Attack (Example)

- <https://pytorch.org/blog/compromised-nightly-dependency>
- PyTorch-nightly Linux packages installed via pip during that time installed a dependency, **torchtriton**, which was compromised on the Python Package Index (PyPI) code repository and ran a malicious binary. This is what is known as a **supply chain attack** and directly affects dependencies for packages that are hosted on public package indices.
- A malicious dependency package (**torchtriton**) that was uploaded to the Python Package Index (PyPI) code repository with the same package name as the one shipped on the [PyTorch nightly package index](#).
- This malicious package was being installed instead of the version from the official repository.
- This malicious package contains code that uploads sensitive data from the machine.

Compromised PyTorch-nightly dependency chain between December 25th and December 30th, 2022.

🔔 If you installed PyTorch-nightly on Linux via pip between December 25, 2022 and December 30, 2022, please uninstall it and torchtriton immediately, and use the latest nightly binaries (newer than Dec 30th 2022).

```
$ pip3 uninstall -y torch torchvision torchaudio torchtriton
$ pip3 cache purge
```

PyTorch-nightly Linux packages installed via pip during that time installed a dependency, torchtriton, which was compromised on the Python Package Index (PyPI) code repository and ran a malicious binary. This is what is known as a supply chain attack and directly affects dependencies for packages that are hosted on public package indices.

NOTE: Users of the PyTorch **stable** packages **are not** affected by this issue.**

Challenges in dependency management

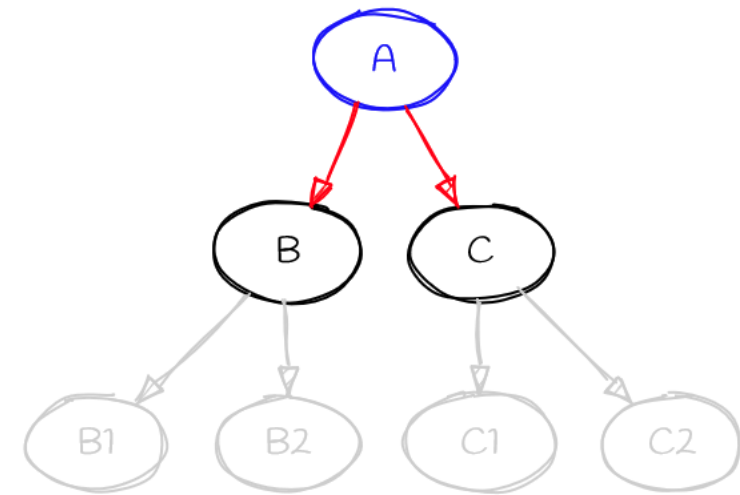
- Management of networks of libraries, packages, and dependencies that we don't control
- Concerns:
 - How to update between versions of external dependencies?
 - How do we describe versions?
 - What changes are allowed/expected?
 - How to decide whether its wise to depend on other org's code?
- Cascade of upgrades

Dependency Resolution

- Given a set of dependencies, find a set of package to install
 - **NP-Hard problem!**
 - The more dependencies you have, the harder and longer it gets to resolve dependencies
- How pip works:
 - **Finder:** Get package data from an “index” (e.g., all pandas versions available)
 - **Resolver:** Considers all “candidates” from finder, and outputs the final list of resolved dependencies

Dependency Resolution

- Resolver algorithm steps (See ResolveLib):
 - **Identify**: identify candidates and requirements
 - **Get Preference**: Selects which requirement to look at “next”
 - Includes a bunch of heuristics
 - **Find matches**: Given a set of constraints, find what candidates exist that satisfy them (uses finder). May use a SAT Solver.
 - **is_satisfied_by**: checks if a candidate satisfies a requirement
 - **Get_dependencies**: get dependency metadata for a candidate.



<https://pip.pypa.io/en/stable/topics/more-dependency-resolution>

<https://pypi.org/project/resolvelib>

Dependency management

- What versions of dependencies should you import?
- When should you upgrade dependency versions?
- SwE@Google book outlines four options:
 - Never upgrade
 - Semantic versioning
 - Bundled distributions
 - "Live at HEAD"

Never upgrade (Static Dependency Model)

- Predictable
 - Avoids failures due to changes outside of your control
- Natural when starting out, or for short-lived projects
 - Compatible with "vendoring"
- What happens when a dependency has a security vulnerability?
- What happens when a new dependency depends on newer versions of old dependencies?

Bundled distributions

- Defer dependency management to the distribution maintainer
 - Responsible for maintaining compatibility while incorporating security updates
- Depend on the bundle and whatever dependency versions it provides
 - Common for commercial applications
- **Example:**
 - Linux distributions: Debian (non-commercial), Fedora Linux (Red Hat), OpenSUSE
 - Android, ChromeOS: Built on Linux Kernel
 - More niche: Raspberry PI OS
- **Distributors:** responsible for finding, patching, and testing a mutually compatible set of versions to include.
- Limitations:
 - Limits (verified) portability
 - Can't leverage latest features

Semantic versioning (SemVer)

- Dependency version numbers obey MAJOR.MINOR.PATCH format
 - Changes to PATCH should be fully compatible (bug fixes, security fixes)
 - Changes to MINOR may add functionality in a backwards-compatible manner
 - Changes to MAJOR indicate API changes (potentially breaking)
- Assumed by many build tools
 - Depend on a specific MAJOR version and a minimum MINOR version
- Challenges
 - Not all dependencies follow this scheme
 - Human maintainers make mistakes
 - **Hyrum's Law**: one person's "bug" is another's "feature"
 - Can be over-constraining (no solution to SAT problem)
 - Heuristics for relaxing some requirements

Minimum Version Selection (MVS)

- SemVer: Chooses the newest possible versions of dependencies that satisfy requirements
- **MVS**: Select the lowest satisfiable version
- **Intuition**: Produce **high-fidelity builds** in which dependencies are as close as possible to what the developer used
- Proposed by Russ Cox for Go: <https://research.swtch.com/vgo-mvs>

"Live at HEAD"

- Dependency management analogue of *trunk-based development*
- **Principles:**
 - Always depend on current stable version of everything
 - Never change anything in a way that is difficult for dependents to adapt
- **Dependency maintainer** responsible for not breaking all users
 - Effectively requires continuous integration for all software in the world (except closed-source dependents)
 - If compatibility cannot be maintained, maintainer will provide an upgrade tool
- **API providers:** Ensure smooth migration; **API consumers:** Provide tests

"Live at HEAD"

- Some of this infrastructure already exists
 - "Rolling" Linux distributions (e.g., Gentoo) integrate tens of thousands of packages continuously
 - Programming languages (e.g. Scala, Rust) proactively test all changes against major libraries/applications
- **Version selection:** What is the latest stable version of everything?

Dependency vulnerabilities

- NPM has a history of dependency-related disasters
 - **left-pad** unpublished
 - Bitcoin theft transitive dependency in **event-stream**
 - Ukraine war "**protestware**" in **node-ipc**
- Why was impact so large?
 - Tools depended on external repository services rather than internal mirror
 - Projects depended on **floating** instead of fixed versions (e.g., >1.5, 5.*)
 - Projects were built "too continuously"
 - Fine-grained dependencies depended upon by many other libraries (cascading)

https://en.wikipedia.org/wiki/Npm_left-pad_incident

<https://www.trendmicro.com/vinfo/us/security/news/cybercrime-and-digital-threats/hacker-infests-node-js-package-to-steal-from-bitcoin-wallets>

<https://orca.security/resources/blog/cve-2022-23812-protestware-malicious-code-node-ipc-npm-package>

Vulnerabilities

- CVE: Common Vulnerabilities and Exposures
 - Common identifier for specific vulnerabilities (not vulnerable systems)
 - CWE: Common Weakness Enumeration (type of vulnerability)
 - May be crosslinked with other databases (e.g., severity, product, weakness category)
 - NIST's National Vulnerability Database (NVD) includes common links and history
 - Common Vulnerability Scoring System (CVSS) standardizes measures of severity
- Example CVE: <https://nvd.nist.gov/vuln/detail/CVE-2025-32955>
- Others: <https://mvnrepository.com/artifact/org.opencontainers/opencontainers-core/16.0>

Reading

- *Software Engineering at Google*, Chapter 21: Dependency Management

Build systems

Build System Objectives

- Automate compilation & linkage of all components
- Rebuild necessary components when things change
- Manage multiple configurations
- Manage external dependencies
- Automate testing
- Automate release actions
 - Strip debugging symbols
 - Minify web assets
 - Generate installers

Also relevant for
interpreted languages

Desirable properties of build system

- Fast:
 - Run a single command to build and get the output binary in a short time (few seconds)
- Correct:
 - Reproducible: Should output same result for any developer/machine for the same input files

What does a build system even do?

- Why something like `javac *.java` is not enough?
- **How to handle:**
 - Building libraries stored in different directories (shared libraries)
 - Code written in different programming languages (dependencies)
 - Third-party jar files (how to store them, version management)
 - Rebuilding part of the codebase after dependency upgrade
 - Target different systems/release builds (build configs)
 - (Implicit dependencies) Managing related artifacts/tasks: documentation, latest library version
- Write a shell script?

Options

- Write your own scripts
 - Lots of redundant effort to provide flexibility and functionality
 - Maintenance cost of bespoke system
- Follow conventions
 - Easy way for new projects to take advantage of build tool features with minimal effort
 - Good IDE support
 - Hard to adapt for large, heterogeneous, legacy projects
 - Difficult to diagnose implicit rules
 - Can lead to bloated dependencies
- Configure a build tool
 - Must learn a complicated tool & configuration syntax
 - But knowledge is transferrable
 - Must maintain build configuration
 - But being explicit is often good, avoids dependency bloat
 - Can accommodate custom procedures
 - Code generation
 - Multiple languages
 - IDE may require additional configuration

Common build tools

- Make [1976]
 - Autoconf
 - CMake
 - Ant + Ivy, Maven, Gradle (Java)
 - sbt (Java, Scala)
 - Pip, setuptools (Python)
 - npm, Bower (Javascript)
 - Cargo (Rust)
 - latexmk (LaTeX)
 - Bazel
- Responsible for constructing **dependency graph**
 - Task-oriented: Targets can execute arbitrary commands
 - Hard to correctly specify when a task does not need to be rerun
 - Hard to parallelize safely
 - Artifact-oriented: Targets must declare inputs, outputs
 - Enables safe caching, parallelization

Task-Based Build Systems

- **Task:** Fundamental Unit of Work
- Tasks can have other tasks as dependencies
- Major build systems: Ant, Maven, Gradle, Grunt, Rake, ...

Make example

- Built-in implicit rules
 - Knows how to compile .cc files to get .o file
 - Uses standard env vars (CXX, CXXFLAGS)
- Compiler provides header dependencies for future use
 - But what if a header with the same name is created elsewhere?
- Does not depend on variable values (static)
- Use .PHONY to declare tasks that don't produce artifacts
- First target is default
- Uses timestamp to detect changes

See [scrambler/c++/Makefile](#)

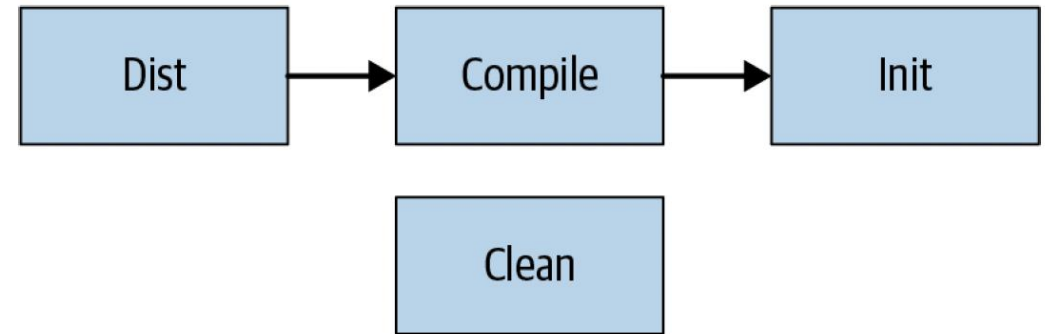
Example: Ant Build File

```
<project name="MyProject" default="dist"
  basedir=".">
<description>
simple example build file
</description>
<!-- set global properties for this build -->
<property name="src" location="src"/>
<property name="build" location="build"/>
<property name="dist" location="dist"/>
<target name="init">
<!-- Create the build directory structure used
  by compile -->
<mkdir dir="${build}"/>
</target>
<target name="compile" depends="init"
description="compile the source">
<!-- Compile the Java code from ${src} into
  ${build} -->
<javac srcdir="${src}" destdir="${build}"/>
</target>
```

```
<target name="dist" depends="compile"
description="generate the
  distribution">
<!-- Create the distribution directory
  -->
<mkdir dir="${dist}/lib"/>
<!-- Put everything in ${build} into
  the MyProject-${DSTAMP}.jar file -->
<jar jarfile="${dist}/lib/MyProject-
  ${DSTAMP}.jar" basedir="${build}"/>
</target>
<target name="clean"
description="clean up">
<!-- Delete the ${build} and ${dist}
  directory trees -->
<delete dir="${build}"/>
<delete dir="${dist}"/>
</target>
</project>
```

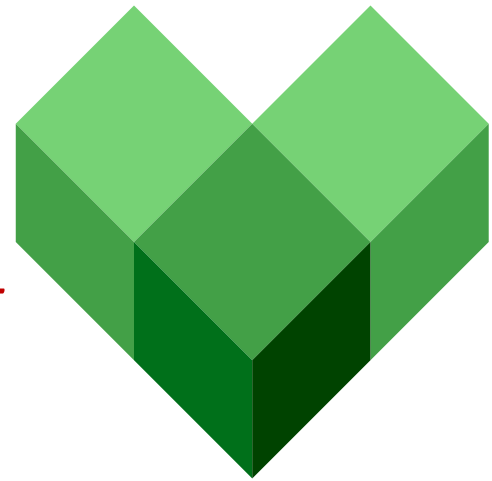
Example: Ant Build File

- “ant [task name]” executes the given task and its dependent tasks
- Advantage:
 - Modularized builds
- Disadvantages:
 - One more file to debug! (tricky)
 - Hard to parallelize
 - Incremental builds are difficult



Artifact-based build system

- **Build:** Tell the system “what” to build instead of “how”
- Implemented in Blaze/Bazel (Google), Pants, Buck
- Build files are declarative: specify set of artifacts to build, their dependencies, some build options (instead of exact steps)
- Bazel has full control over “how” build is run
- *(Stronger) correctness guarantee while being more efficient*



Example Bazel BUILD file

```
java_binary(  
  name = "MyBinary",  
  srcs = ["MyBinary.java"],  
  deps = [ ":mylib", ],)  
java_library(  
  name = "mylib",  
  srcs = ["MyLibrary.java", "MyHelper.java"],  
  visibility =  
    ["//java/com/example/myproduct:__subpackage  
    s__"],  
  deps = [  
    "//java/com/example/common",  
    "//java/com/example/myproduct/otherlib",  
    "@com_google_common_guava_guava//jar",  
  ], )
```

- **Targets/Artifacts:**
java_binary, java_library
- **Workspace:** Source hierarchy for artifacts

Bazel BUILD Steps

```
bazel build :MyBinary
```

- Parse all build files and create graph of artifacts and dependencies
- Determine transitive dependencies of MyBinary
- Build each dependency (in order).
 - Start with artifacts with no dependencies.
 - Keep track of artifacts that need dependencies to be built
 - Build a target **as soon as** its dependencies are built
- Build final MyBinary executable binary

Bazel Advantages/Differences

- Parallelization:
 - Targets that only require java compiler (vs custom script)
- Reuse/caching:
 - If MyBinary.java changes, it will rebuild **MyBinary** but reuse **mylib**
 - If a source file for **//java/com/example/common** changes, Bazel knows to rebuild that library, **mylib**, and **MyBinary**, but reuse **//java/com/example/myproduct/otherlib**

PollEv.com/cs5150sp26

Which of the following best describes an advantage of Bazel over traditional build systems like Make or custom Java build scripts?

- A.** Bazel always rebuilds the entire project to ensure consistency.
- B.** Bazel parallelizes tasks using custom shell scripts instead of native compilers.
- C.** Bazel tracks fine-grained dependencies, enabling it to rebuild only what's necessary and reuse cached outputs.
- D.** Bazel relies on environment variables for dependency tracking and compilation.

Other Bazel Features

- Tools as dependencies, **toolchains** (platform-specific tool usage)
- Custom user-defined **actions**: specify inputs, outputs, and steps
- **Sandboxing**: isolating filesystem for each action
- Remote caching
- **Distributed build**: Remote build
- Making remote/external dependencies deterministic
 - Manifest file: Create **cryptographic hash** for each ext dependency, only redownload when hash changes, build fails if hash changes
 - What can go wrong?

Dependency Management @ Google

Scaling to **Billions** of Lines of Code

- Strict Transitive Dependency:
 - A cannot use a symbol for C without declaring direct dependency
- External Dependencies: Uses semantic versioning
 - One-Version Rule (eliminates diamond dependency problem)
- Transitive External Dependencies:
 - Bazel does not allow automatic download of such dependencies
- Shared cache for external artifacts that require building
- Security/Reliability: Mirroring servers, Vendoring

