# Reliable Benchmarking of Software Verification in the Cloud

Philipp Wendler

**SoSy-Lab**
Software Systems

UNIVERSITÄT
PASSAU

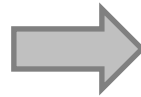*Fakultät für Informatik und Mathematik*

# Software Verification

Goal: Build an automatic software verifier

C program

```
int main() {
 int x = 10;
 int y = 3;
 int z = x+y;
 assert(z > 0);
}
```

Specification

Software
Verifier

**SAFE**
i.e., assertions
cannot be violated

**UNSAFE**
i.e., bug
in program

Our tool: CPAchecker
http://cpachecker.sosy-lab.org
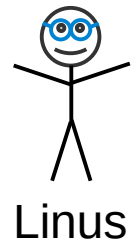
CPA ✓

Linux Driver Verification

Photo: (c) Sean Bonner, 2013

# Weekly Verification Workload

## Time to verify...

- 1300 commits
- 4 drivers affected (approximately)
- 100 (safety) properties
- 12 seconds per verification task (on average)

Way more than 100 properties should be checked!

Linus

**1700 hours,** or **72 days** of CPU time

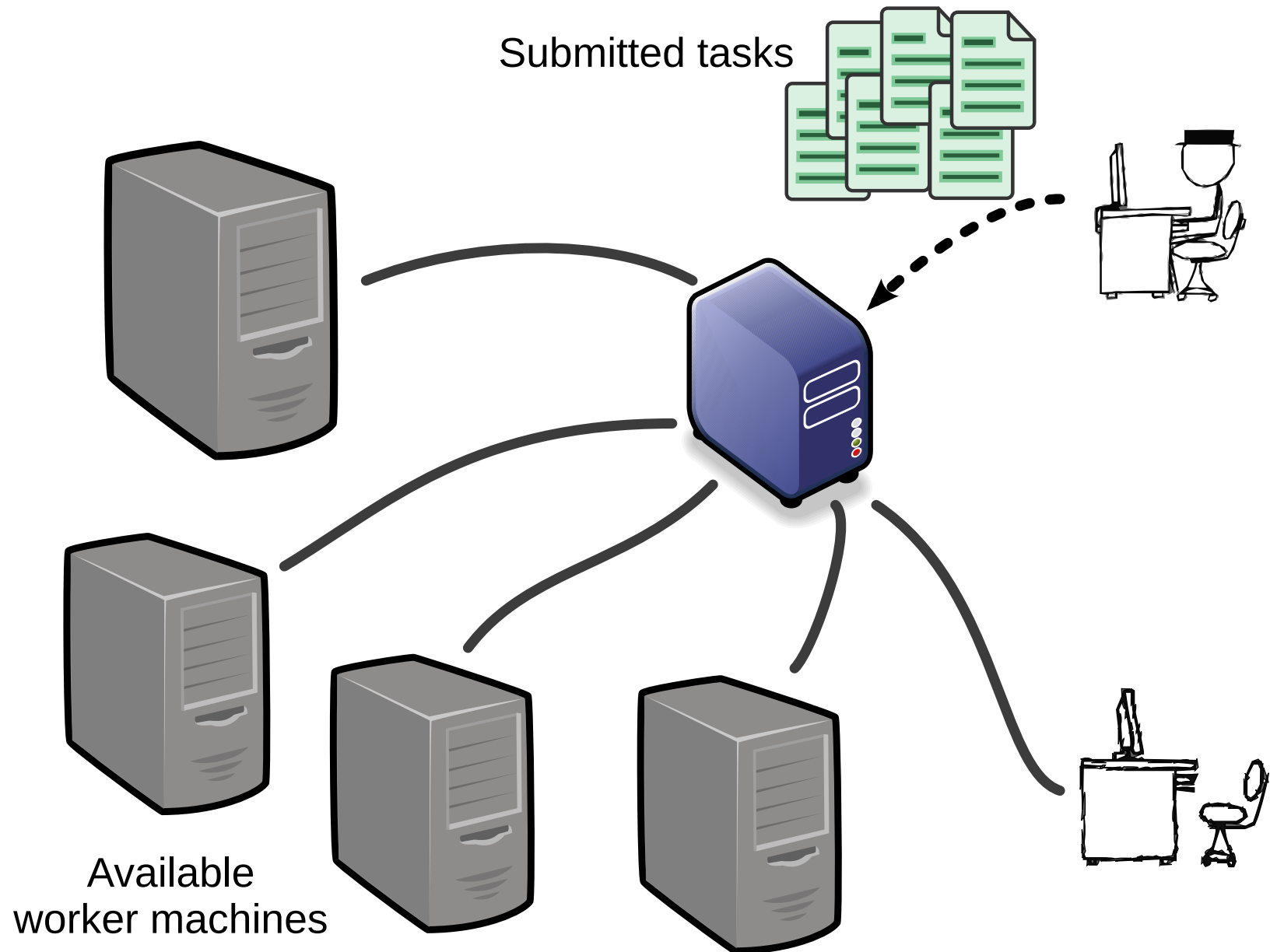# Software Verification in the Cloud

Verification of large number of tasks

Reliable benchmarking:
Accurate and reproducible results

# VerifierCloud: Design

- Distribution system for verification tasks

- Task submission, e.g., via web frontend

- Creates necessary execution environment

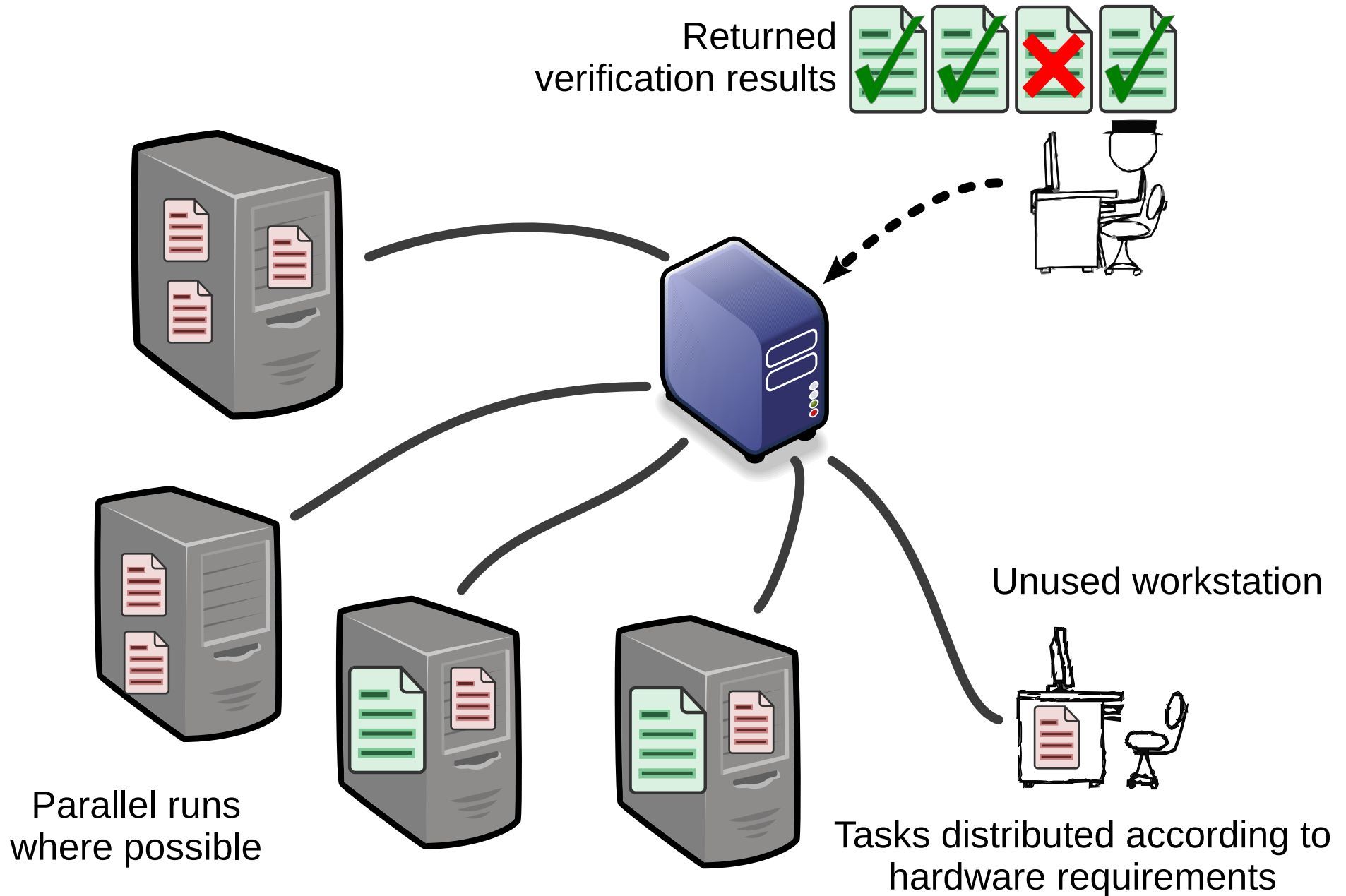- Supports arbitrary worker machines
(no shared file system necessary)

# VerifierCloud: Architecture

Submitted tasks

Available
worker machines

# VerifierCloud: Machine Allocation

- Partitions hardware resources
  and executes multiple runs
  in parallel on each machine

- Worker machines can be used dynamically
  when not needed by other users

# VerifierCloud: Architecture



Returned verification results

Parallel runs where possible

Unused workstation

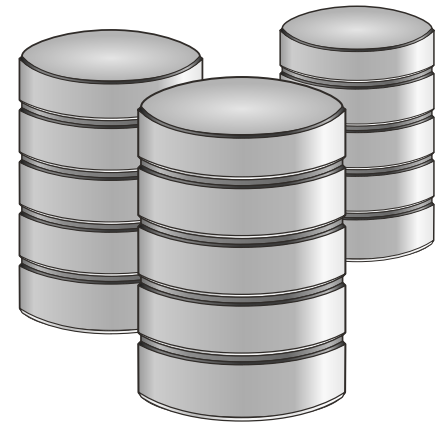Tasks distributed according to hardware requirements

# VerifierCloud: Results

- Successful use in development of our verifier CPAchecker

    - Multiple groups, ca. 40 developers

    - Significantly speeds up implement-test-roundtrip time

    - Up to 800 runs executing in parallel

    - 1 000 000 runs per week

- Used in teaching (Passau, Paderborn, Hamburg)

    - Students can learn verification without installation

    - Resources for experimenting with own implementations

# VerifierCloud: Future

- Reuse of results
  - Similar verification tasks can benefit from reusing intermediate results
  - Example: new revision of same driver
- Store (intermediate) results in database?
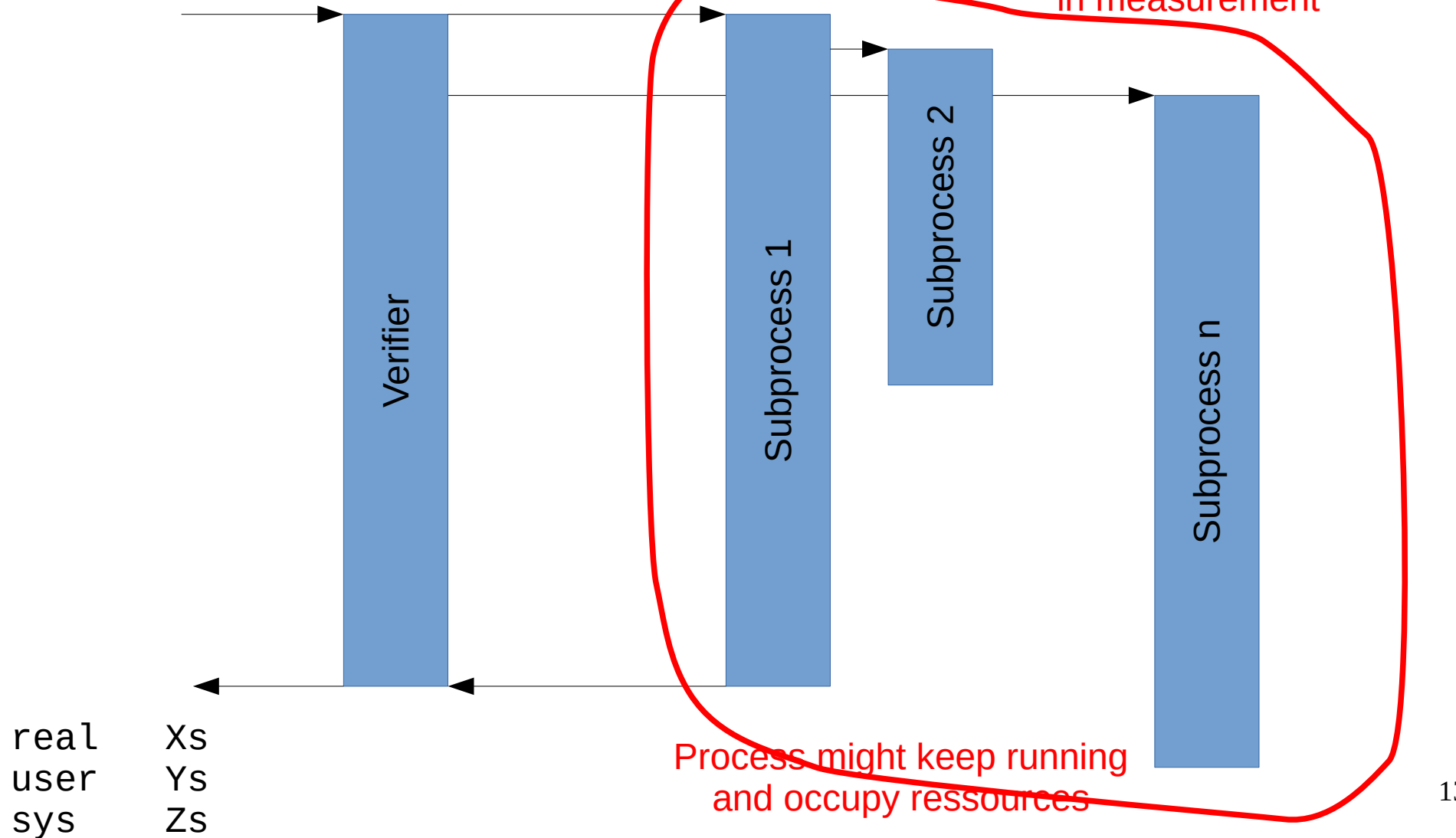- Automatically use stored information for new tasks?

# Reliable Benchmarking

- Shared Machines

- Arbitrary Tools

  - Non-interactive

  - CPU-/Memory bound (no I/O)

- Measure and limit CPU time

- Measure and limit memory
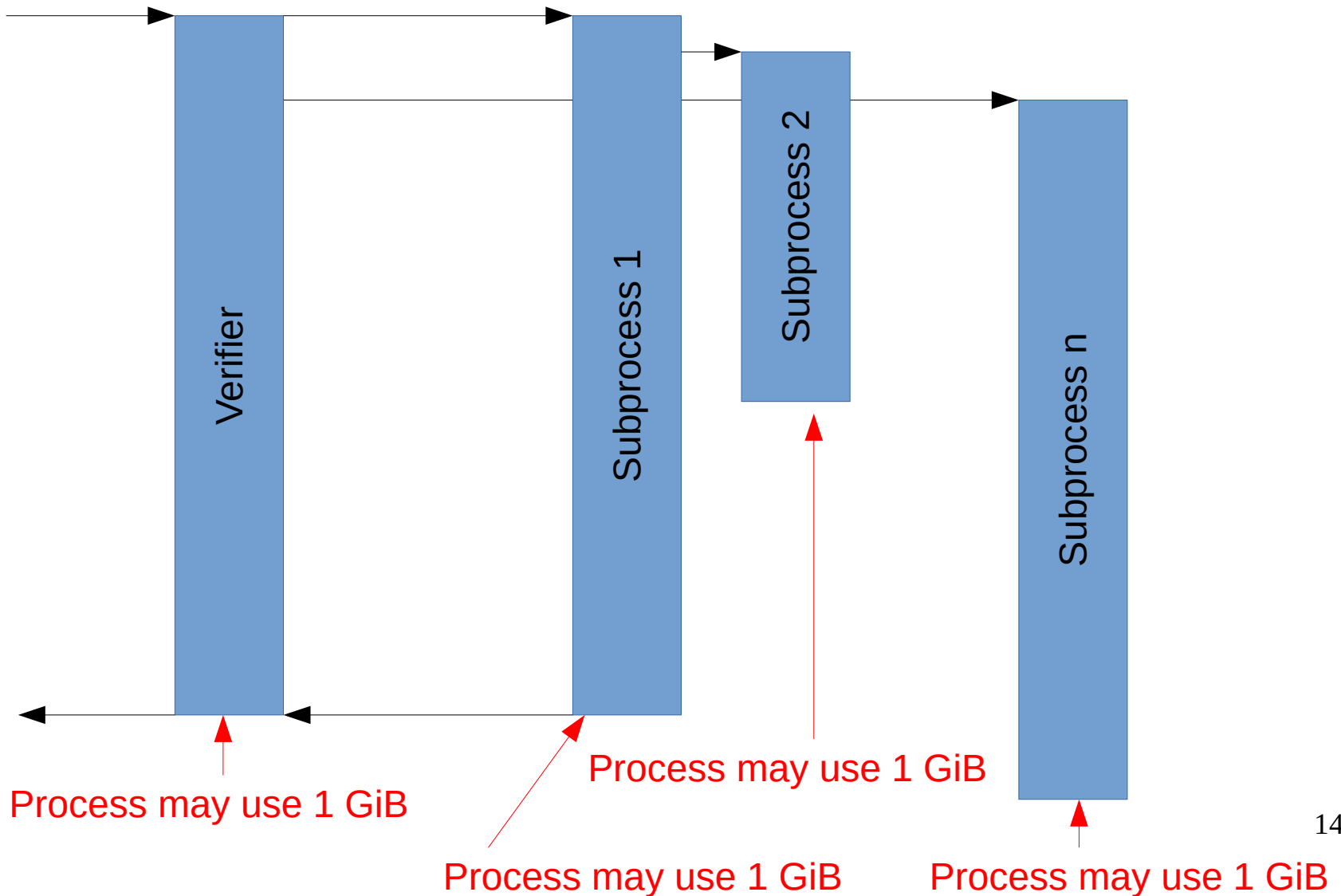
- Different Hardware Architectures
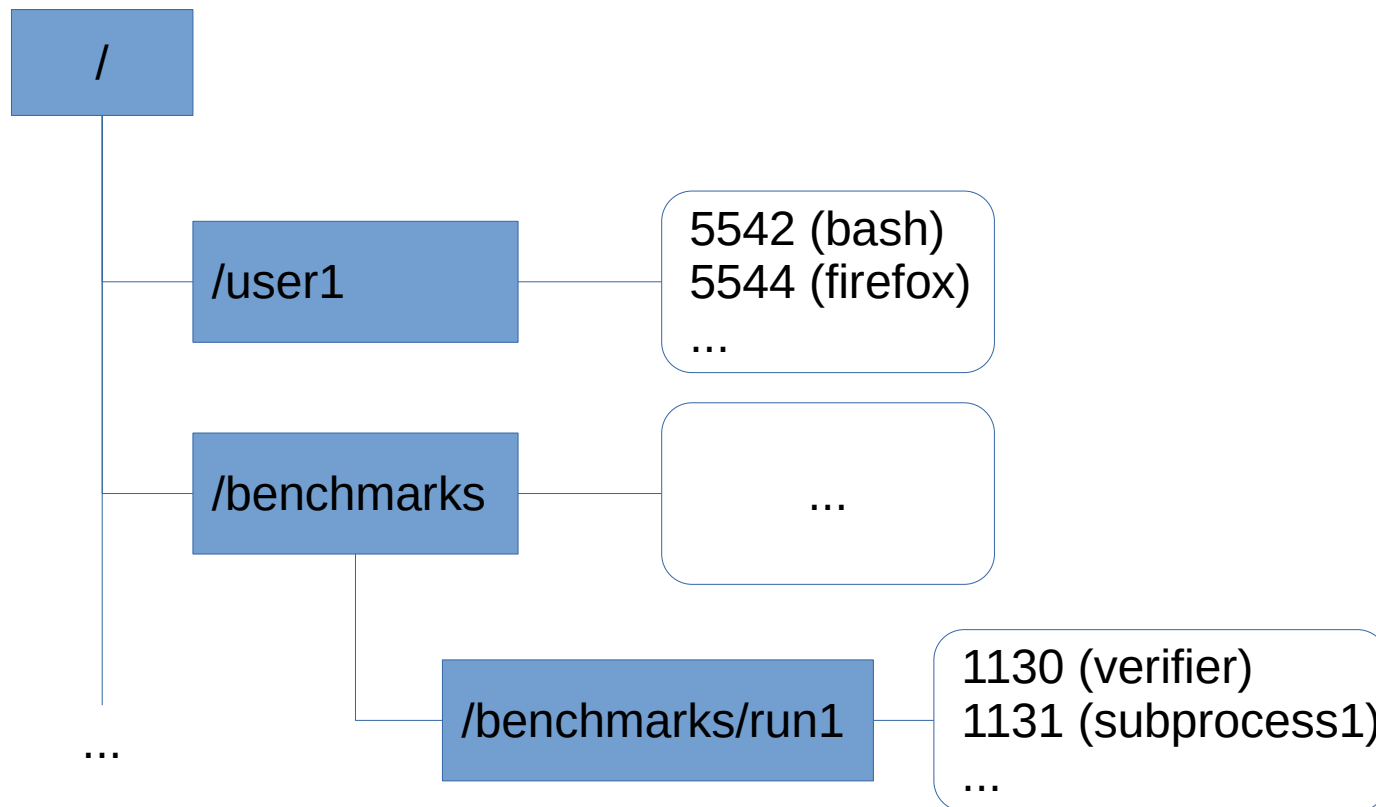
# Measuring CPU time with „`time`“

`~$ time verifier`

Verifier

Subprocess 1

Subprocess 2

Subprocess n

CPU Time may not be included in measurement

Process might keep running and occupy ressources

```
real    Xs
user    Ys
sys     Zs
```

# Limiting memory with „ulimit"

```
~$ ulimit -v 1048576 # 1 GiB
~$ verifier
```



Verifier

Subprocess 1

Subprocess 2

Subprocess n

Process may use 1 GiB

Process may use 1 GiB

Process may use 1 GiB

Process may use 1 GiB

# Cgroups

- Linux kernel „control groups"
- Hierarchical tree of sets of processes

```
/
├── /user1 ─── 5542 (bash)
│              5544 (firefox)
│              ...
├── /benchmarks ─── ...
│       └── /benchmarks/run1 ─── 1130 (verifier)
│                                1131 (subprocess1)
│                                ...
...
```

# Cgroups

- Reliable tracking of spawned processes

- Resource limits and measurements per cgroup
  - CPU time
  - Memory
  - I/O etc.

Only solution on Linux
for race-free handling of multiple processes!

# Hardware Architectures

- Hyper Threading

  Multiple threads sharing execution units

- Non-Uniform Memory Access

  Memory regions have different performance depending on current CPU core

- And more (caches, …)

- Can lead to non-deterministic performance

# BenchExec

- A Framework for Reliable Benchmarking and Resource Measurement

- Based on cgroups

- Handles multiple processes

- Allocates hardware resources appropriately

- Used in International Competition on Software Verification (SV-COMP)

  22 tools this year

# BenchExec

- Open source: Apache 2.0 License
- https://github.com/dbeyer/benchexec
- Paper under submission
- Extensible
  - Arbitrary tools
  - Not only for software verification