



université  
PARIS-SACLAY



# Analyzing compile-time overhead

Hadrien Grasland

2022-07-05



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 101004761.



# First, a little story

- In ACTS, we tried to follow **Modern C++ Best Practices™**
  - Avoid vtables like pest, use templates all the way down
  - Put most code in headers because templates + inlining
  - Use fancy metaprograms to reduce runtime work

# First, a little story

- In ACTS, we tried to follow **Modern C++ Best Practices™**
  - Avoid vtables like pest, use templates all the way down
  - Put most code in headers because templates + inlining
  - Use fancy metaprograms to reduce runtime work
- It worked ok for years, then **weird CI crashes** crept in
  - Job abruptly fails, not always in the same place, no error log
  - Reducing concurrency for reproductibility, it goes away
  - ...until someday it doesn't and the crashes are back at -j1



# Investigation begins

- « Luckily », several developers got **similar local crashes**
  - Easier to figure out they're running out of RAM
  - Often decide to live with slower builds by tuning down -j
  - Friendly colleague suggests **downloading more RAM**

# Investigation begins

- « Luckily », several developers got **similar local crashes**
  - Easier to figure out they're running out of RAM
  - Often decide to live with slower builds by tuning down -j
  - Friendly colleague suggests **downloading more RAM**
- ...but none of these workarounds will work for Github CI VMs !
  - Need to admit that **our codebase is the problem**
  - Problem is, there are about 500k lines of code in it
  - How to efficiently figure out which of these are the culprits ?

# **Build profiling howto**

# Step 1 : Set a goal

- Minimal sanity : **Don't crash CI**
  - If you use Github CI, you have 2 cores & 7 GB
  - So even -j1 build must use <6 GB (keeping 1GB for OS)
  - If you want -j2 builds, get below 3 GB/process

# Step 1 : Set a goal

- Minimal sanity : **Don't crash CI**
  - If you use Github CI, you have 2 cores & 7 GB
  - So even -j1 build must use <6 GB (keeping 1GB for OS)
  - If you want -j2 builds, get below 3 GB/process
- Optimal sanity : **Use dev machines efficiently**
  - Classic mid-end laptop has 2 GB RAM per (logical) core
  - So for 4 threads / 8 GB, try to get below 1,5 GB/process\*

\* The typical dev computer has more background tasks → 2GB for OS+background is safer



## Step 2 : Pick pathological files

- There is no reliable\* tool to analyze your full build
  - Best to focus on individual compilation units (cpp files)
- How to find bloated compilation units ?
  - Easy but tedious : Watch system monitor during build
  - Automated : Use **cmakeperf** (Paul's part)

\* <https://github.com/aras-p/ClangBuildAnalyzer> is nice but does not scale...

# Step 3 : Find a proxy for RAM usage

- Compilers won't tell you anything about RAM consumption
  - You don't want to run e.g. heaptrack on those beasts\*
- However, **what consumes RAM, most likely takes time**
  - And we do have time profiles, via **clang's -ftime-trace**
  - Well, not quite a time profile, but an activity history
  - However, the information we need is hidden in there

\* Learning that memory is consumed while instantiating *some* templated code is not useful



# Step 5 : Profit !

Hierarchical profile:



# Step 5 : Profit !

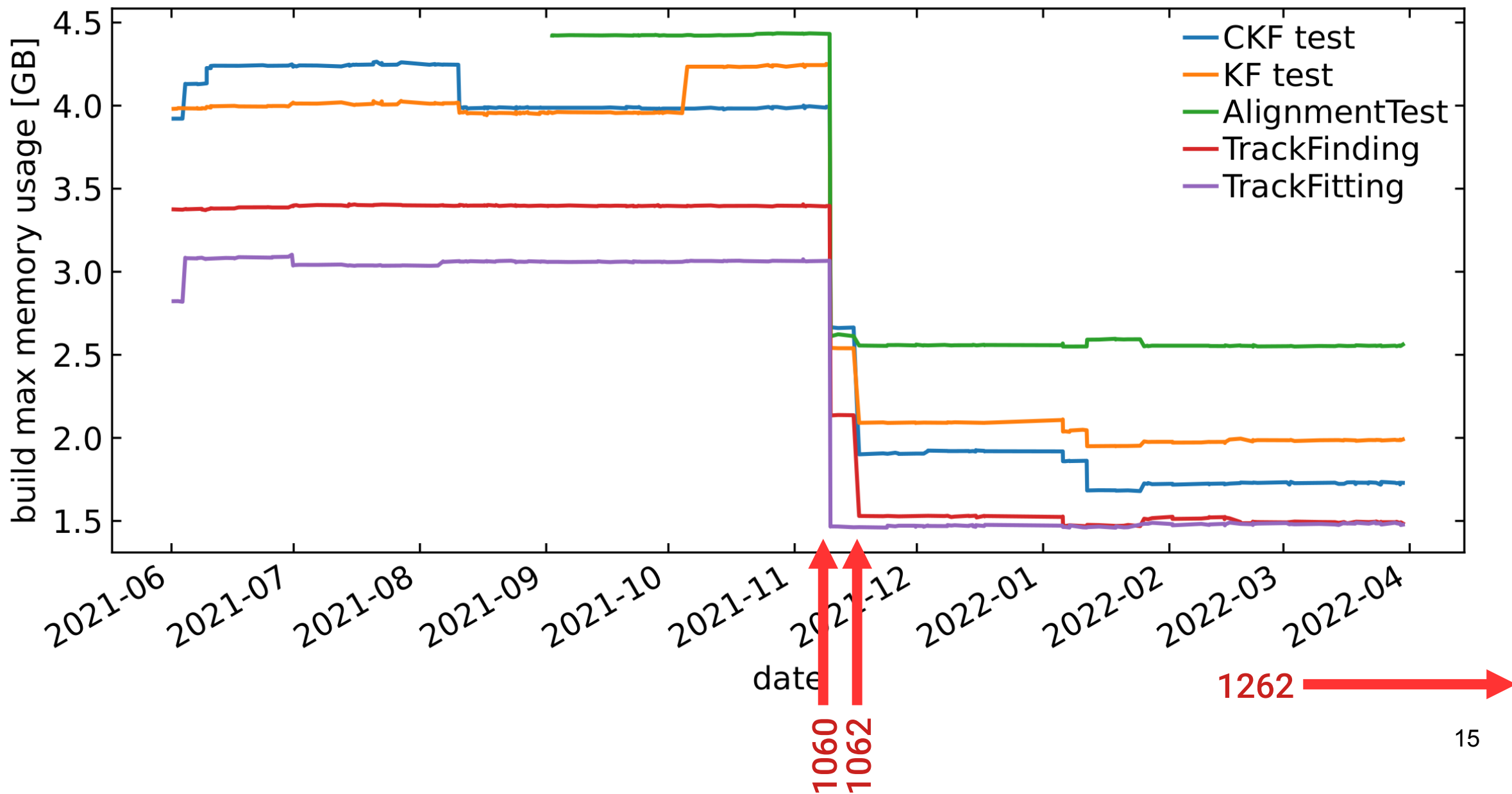
Hierarchical profile:



...and this also has a WIP perf report style interactive viz :)

# **Application to ACTS**

# Let's look at some PRs

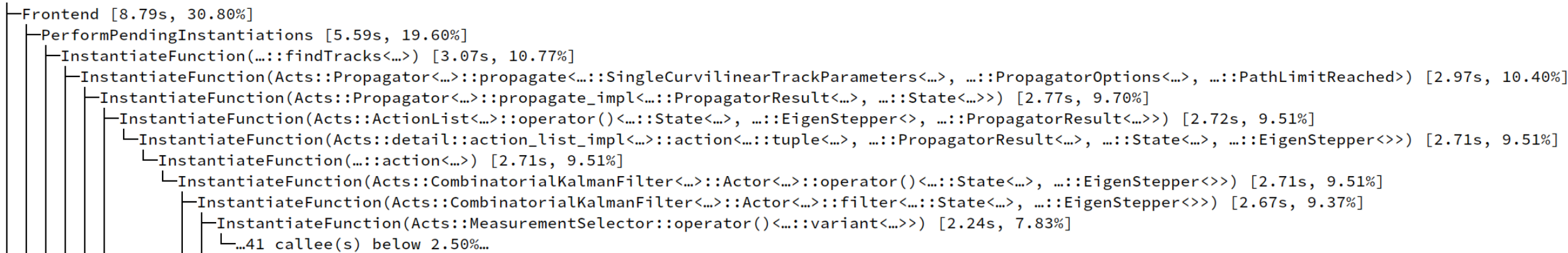


# Before Acts#1060 (SourceLink rework)





# After Acts#1060 (SourceLink rework)



- Using the same temporal display cutoff as before (~800ms)
  - Frontend time was reduced by 5,71s (-39 %)
  - Max-RSS was reduced by 1,3 GB (-32 %)
- Why ?
  - Mostly loss of CKF::Actor::addSourceLinkState (-30%)
  - As Paul will explain, it could be moved to a dedicated CU

# Before Acts#1062 (Delegate @ CKF)



# After Acts#1062 (Delegate @ CKF)

```
└─Frontend [6.56s, 32.18%]  
  └─PerformPendingInstantiations [3.29s, 16.13%]  
    └─...294 callee(s) below 3.70%...
```

- Again, used similar temporal cutoff for fair comparison
  - But to find back CKF::findTrack, we'll need a smaller cutoff

```
└─Frontend [6.56s, 32.18%]  
  └─PerformPendingInstantiations [3.29s, 16.13%]  
    └─InstantiateFunction(Acts::Test::createMeasurements<...::Propagator<...>, ...::SingleCurvilinearTrackParameters<...>>) [600.69ms, 2.95%]  
      └─InstantiateFunction(Acts::Propagator<...>::propagate<...::SingleCurvilinearTrackParameters<...>, ...::PropagatorOptions<...>, ...::PathLimitReached>) [564.31ms, 2.77%]  
        └─...15 callee(s) below 2.38%...  
        └─...7 callee(s) below 2.38%...  
    └─InstantiateFunction(Acts::CombinatorialKalmanFilter<...>::findTracks<...::TestContainerAccessor<...>, ...::vector<...>, ...::SingleBoundTrackParameters<...>>) [486.94ms, 2.39%]  
      └─...17 callee(s) below 2.38%...
```

- -87 % decrease → no longer the frontend bottleneck. How ?
  - MeasurementSelector could be moved to a different CU
  - The Eigen pyramid of doom followed it

# The Future : Acts#1262

- Acts being integrated in ATLAS → Need to integrate w/ xAOD
  - Need more genericity in trajectory storage → sensitive !

	file	max_rss_main	max_rss_feat	relative
10	CombinatorialKalmanFilterTests.cpp	1087.35	1230.05	113.123
11	GsfTests.cpp	1642.13	1686.95	102.729
12	KalmanFitterTests.cpp	1232.11	1282.43	104.085
13	GainMatrixUpdaterTests.cpp	383.144	437.797	114.264
14	GainMatrixSmootherTests.cpp	329.216	413.655	125.649

- Paul managed to get it down to a ~100MB worst-case cost
  - Where does this come from, and could we do better ?

# GainMatrixSmootherTests before

Hierarchical profile:

ExecuteCompiler [3.77s, 100.00%]

└─Frontend [2.42s, 64.24%]

└─Source(/opt/spack/opt/spack/.../gcc-12.1.0/boost-1.79.0-pu25xp7wrflbj2xwdb3dcoftqtnde4qg/include/boost/test/unit\_test.hpp) [775.90ms, 20.58%]

└─Source(/opt/spack/opt/spack/.../gcc-12.1.0/boost-1.79.0-pu25xp7wrflbj2xwdb3dcoftqtnde4qg/include/boost/test/test\_tools.hpp) [711.90ms, 18.89%]

└─┬...9 callee(s) below 15.00%...

└─┬...1 callee(s) below 15.00%...

└─PerformPendingInstantiations [750.87ms, 19.92%]

└─┬...115 callee(s) below 15.00%...

└─┬...10 callee(s) below 15.00%...

└─Backend [1.11s, 29.57%]

└─Optimizer [789.54ms, 20.95%]

└─ModuleInlinerWrapperPass(/mnt/acts/Tests/UnitTests/Core/TrackFitting/GainMatrixSmootherTests.cpp) [636.96ms, 16.90%]

└─ModuleToPostOrderCGSCCPassAdaptor(/mnt/acts/Tests/UnitTests/Core/TrackFitting/GainMatrixSmootherTests.cpp) [632.76ms, 16.79%]

└─┬...238 callee(s) below 15.00%...

└─┬...2 callee(s) below 15.00%...

└─┬...13 callee(s) below 15.00%...

└─┬...1 callee(s) below 15.00%...

└─┬...1 callee(s) below 15.00%...

# GainMatrixSmootherTests after

Hi Eigen !

Hierarchical profile:

ExecuteCompiler [5.51s, 100.00%]

Frontend [3.15s, 57.07%]

PerformPendingInstantiations [1.43s, 25.89%]





# What does this give us ?

- Now we have a backtrace of template instantiations
  - Walk through it from the bottom
  - Look for dynamic dispatch & outlining opportunities
  - Check runtime impact wrt expected call frequency



# Lessons learned

- Compilation can have **bottlenecks** just like execution
  - If you can find those, optimizing is much easier
  - Compile time isn't a bad proxy for compile RSS



# Lessons learned

- Compilation can have **bottlenecks** just like execution
  - If you can find those, optimizing is much easier
  - Compile time isn't a bad proxy for compile RSS
- Mind the run-time vs compile-time **tradeoffs**
  - Outlining or vtable *at the right layer* is often good enough
  - If you can't split code in CUs, you'll die by combinatorics
  - Eigen-style cleverness does not scale → At least isolate it\*

\* Or, better yet, remove it after checking what you get for the price.

# Lessons learned

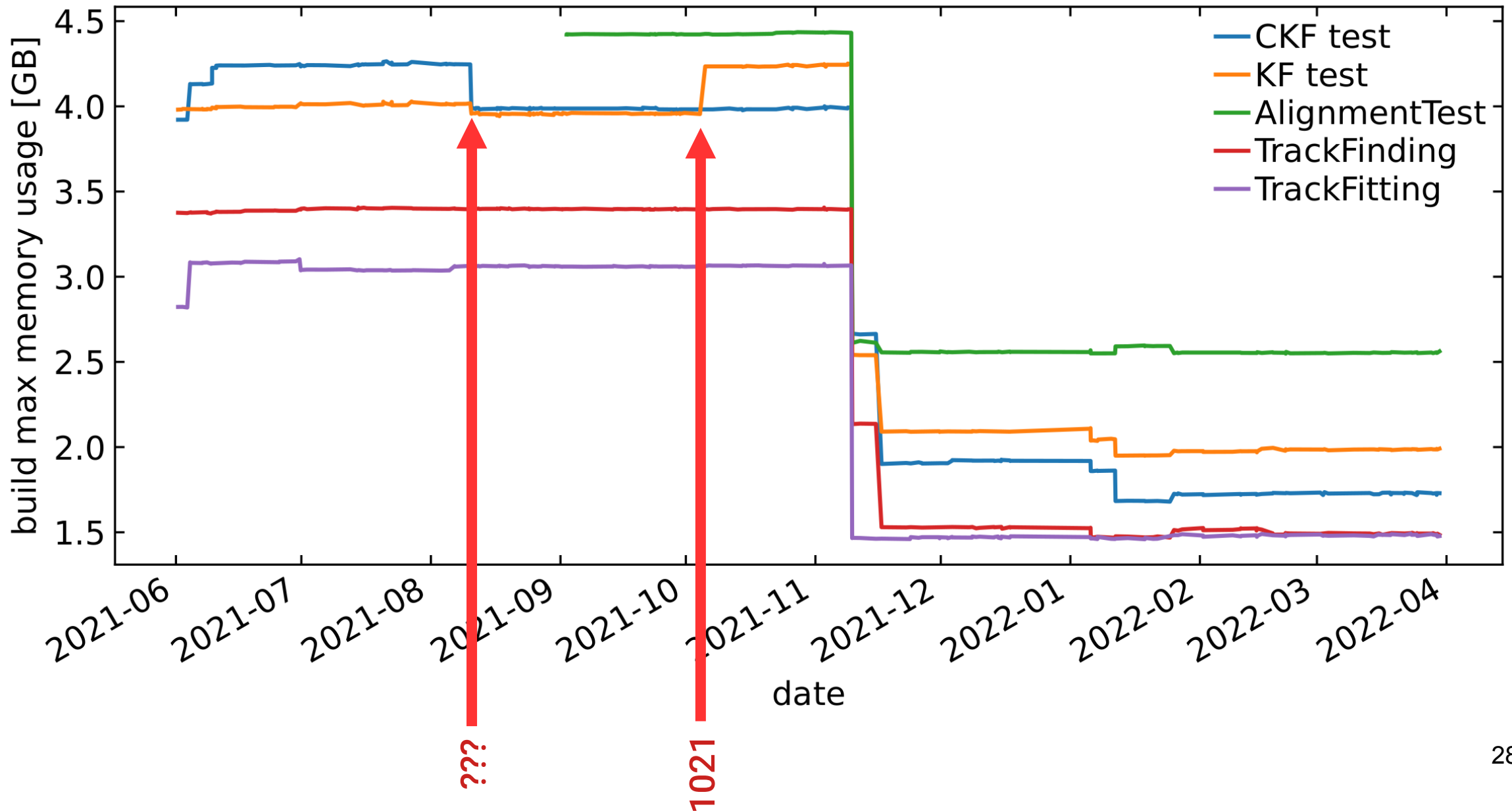
- Compilation can have **bottlenecks** just like execution
  - If you can find those, optimizing is much easier
  - Compile time isn't a bad proxy for compile RSS
- Mind the run-time vs compile-time **tradeoffs**
  - Outlining or vtable *at the right layer* is often good enough
  - If you can't split code in CUs, you'll die by combinatorics
  - Eigen-style cleverness does not scale → At least isolate it\*
- Compilation performance can and should be **monitored**

\* Or, better yet, remove it after checking what you get for the price.

**Paul's turn !**

(Or backup)

# Want some more ?



# Acts#1021 regression

- This PR adds a reverse filtering option to the KalmanFitter
  - Implemented as yet another template parameter
  - Provided with matching unit tests
- It results in a +200MB build-RSS bump in tests
  - Is this just due to extra testing code, or will it affect users ?

# Before Acts#1021

Hierarchical profile:

ExecuteCompiler [44.34s, 100.00%]

└Backend [25.16s, 56.75%]

└└Optimizer [18.05s, 40.71%]

└└└ModuleInlinerWrapperPass(/mnt/acts/Tests/UnitTests/Core/TrackFitting/KalmanFitterTests.cpp) [13.49s, 30.42%]

└└└└ModuleToPostOrderCGSCCPassAdaptor(/mnt/acts/Tests/UnitTests/Core/TrackFitting/KalmanFitterTests.cpp) [13.33s, 30.07%]

└└└└└...4351 callee(s) below 15.00%...

└└└└└...2 callee(s) below 15.00%...

└└└└└...18 callee(s) below 15.00%...

└└CodeGenPasses [7.10s, 16.00%]

└└└OptModule(/mnt/acts/Tests/UnitTests/Core/TrackFitting/KalmanFitterTests.cpp) [7.06s, 15.92%]

└└└└...487 callee(s) below 15.00%...

└Frontend [12.28s, 27.68%]

└└PerformPendingInstantiations [9.05s, 20.41%]

└└└...254 callee(s) below 15.00%...

└└└...37 callee(s) below 15.00%...

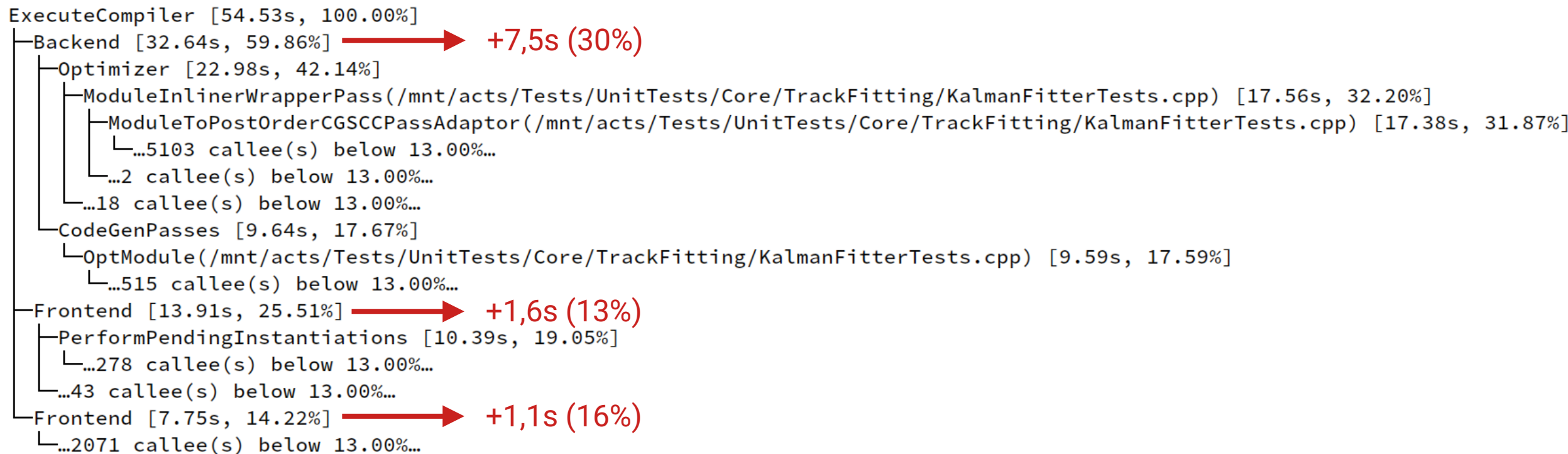
└Frontend [6.70s, 15.10%]

└└...1863 callee(s) below 15.00%...

- This time we're gonna need a higher level view...

# After Acts#1021

Hierarchical profile:



- This time we're gonna need a higher level view...
  - Mostly adds compiler backend work
  - Optimizer hard to analyze (tiny passes) → check CodeGen

# CodeGen before Acts#1021

```
└─CodeGenPasses [7.10s, 16.00%]
  └─OptModule(/mnt/acts/Tests/UnitTests/Core/TrackFitting/KalmanFitterTests.cpp) [7.06s, 15.92%]
    └─OptFunction(Acts::KalmanFitter<...>::Actor<...>::operator()<...::State<...>, ...::EigenStepper<...>> (...) const) [398.68ms, 0.90%]
      └─RunPass(X86 DAG->DAG Instruction Selection) [96.05ms, 0.22%]
      └─RunPass(Live DEBUG_VALUE analysis) [89.04ms, 0.20%]
      └─...56 callee(s) below 0.18%...
    └─OptFunction(Acts::Propagator<...>::propagate<...::SingleCurvilinearTrackParameters<...>, ...::PropagatorOptions<...>, ...::PathLimitReached> (...) const) [194.46ms, 0.44%]
      └─...37 callee(s) below 0.18%...
    └─OptFunction(Acts::Navigator::target<...::State<...>, ...::EigenStepper<...>> (...) const) [180.92ms, 0.41%]
      └─...40 callee(s) below 0.18%...
    └─OptFunction(...::operator() ...) [132.25ms, 0.30%]
      └─...31 callee(s) below 0.18%...
    └─OptFunction(Acts::Navigator::status<...::State<...>, ...::EigenStepper<...>> (...) const) [127.75ms, 0.29%]
      └─...30 callee(s) below 0.18%...
    └─OptFunction(...::test_method ...) [87.25ms, 0.20%]
      └─...25 callee(s) below 0.18%...
    └─OptFunction(Acts::KalmanFitter<...>::Actor<...>::filter<...::State<...>, ...::EigenStepper<...>> (...) const) [80.46ms, 0.18%]
      └─...20 callee(s) below 0.18%...
```



# CodeGen after Acts#1021

```
└─CodeGenPasses [9.64s, 17.67%]
  └─OptModule(/mnt/acts/Tests/UnitTests/Core/TrackFitting/KalmanFitterTests.cpp) [9.59s, 17.59%]
    └─OptFunction(Acts::KalmanFitter<...>::Actor<...>::operator()<...::State<...>, ...::EigenStepper<...>> (...) const) [420.68ms, 0.77%]
      └─RunPass(X86 DAG->DAG Instruction Selection) [99.78ms, 0.18%]
      └─RunPass(Live DEBUG_VALUE analysis) [92.95ms, 0.17%]
      └─...66 callee(s) below 0.15%...
    └─OptFunction(Acts::KalmanFitter<...>::Actor<...>::operator()<...::State<...>, ...::EigenStepper<...>> (...) const) [419.40ms, 0.77%]
      └─RunPass(X86 DAG->DAG Instruction Selection) [104.00ms, 0.19%]
      └─RunPass(Live DEBUG_VALUE analysis) [92.75ms, 0.17%]
      └─...66 callee(s) below 0.15%...
    └─OptFunction(Acts::Propagator<...>::propagate<...::SingleCurvilinearTrackParameters<...>, ...::PropagatorOptions<...>, ...::PathLimitReached> (...) const) [233.58ms, 0.43%]
      └─...41 callee(s) below 0.15%...
    └─OptFunction(Acts::Propagator<...>::propagate<...::SingleCurvilinearTrackParameters<...>, ...::PropagatorOptions<...>, ...::PathLimitReached> (...) const) [223.69ms, 0.41%]
      └─...39 callee(s) below 0.15%...
    └─OptFunction(Acts::Navigator::target<...::State<...>, ...::EigenStepper<...>> (...) const) [212.28ms, 0.39%]
      └─...47 callee(s) below 0.15%...
    └─OptFunction(Acts::Navigator::target<...::State<...>, ...::EigenStepper<...>> (...) const) [211.01ms, 0.39%]
      └─...45 callee(s) below 0.15%...
    └─OptFunction(Acts::Navigator::status<...::State<...>, ...::EigenStepper<...>> (...) const) [148.79ms, 0.27%]
      └─...35 callee(s) below 0.15%...
    └─OptFunction(Acts::Navigator::status<...::State<...>, ...::EigenStepper<...>> (...) const) [146.13ms, 0.27%]
      └─...36 callee(s) below 0.15%...
    └─OptFunction(...::operator() ...) [143.24ms, 0.26%]
      └─...33 callee(s) below 0.15%...
    └─OptFunction(...::test_method ...) [100.35ms, 0.18%]
      └─...25 callee(s) below 0.15%...
    └─OptFunction(Acts::KalmanFitter<...>::Actor<...>::filter<...::State<...>, ...::EigenStepper<...>> (...) const) [97.56ms, 0.18%]
      └─...25 callee(s) below 0.15%...
```

- More complex test → More template instantiations
  - Some of the associated work is duplicated
  - Should not affect production (only 1 KF instantiated there)



# The Great August 2021 Mystery

- Monitoring observed a 200-300MB RSS reduction then
- Does not reproduce locally with clang 14 or gcc 12.1
- Timestamp matches [PR 905](#), which seems largely unrelated
- Maybe a change in CI environment ?