

Application parallelization for multi-core Android devices

Klaas van Gend
klaas@vectorfabrics.com

LinuxCon / ELC Europe
Nov. 6, 2012, Barcelona, Spain



Let me introduce myself + Vector Fabrics



Name: Klaas van Gend

Founded: 1973

History / Claim to fame:

- Started programming C at 14
- First experience with Linux in 1993
- Co-founder of ELC-Europe

Name: Vector Fabrics

Founded: 2007

Claim to fame:

2012: Pareon tool •
analyse & optimize for multicore •



Multicore ARM and Android conquer the world



Google Nexus 10
2-core Samsung A15



Asus Transformer Prime
"4"-core Nvidia Tegra3



HTC J Butterfly
4-core Qualcomm



Samsung Galaxy SIII
4-core Samsung



Sony Xperia P
2-core ST-Ericsson



Huawei Honor2
4-core Huawei

Multi-core usage in Mobile

- 2 core processors:
Assume the OS has **multiple processes** and/or kernel threads to occupy the two cores. **Easy!**
- 4 core processors (and beyond):
Requires **multi-threaded applications** **Hard!**
 - To obtain sufficient concurrent workload
 - To obtain top user experience

Who makes such applications??

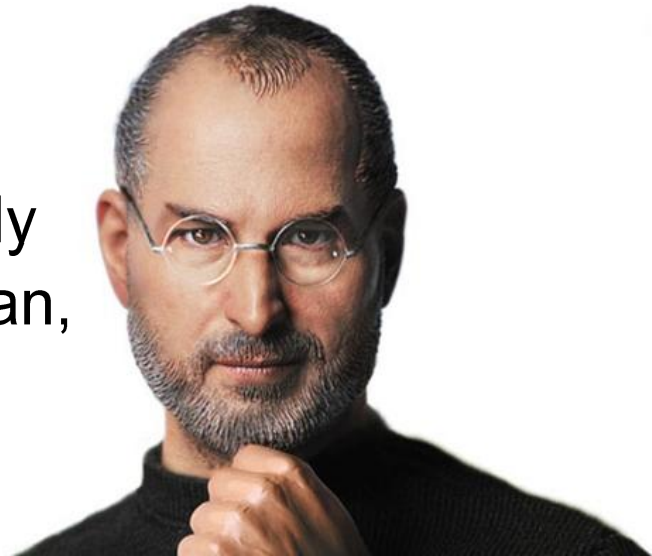
Creating parallel programs is hard...

Herb Sutter, chair of the ISO C++ standards committee,
Microsoft:

“Everybody who learns concurrency thinks they understand it, ends up finding mysterious races they thought weren’t possible, and discovers that they didn’t actually understand it yet after all”

Steve Jobs, Apple:

“The way the processor industry is going, is to add more and more cores, but nobody knows how to program those things. I mean, two yeah; four not really; eight, forget it.”



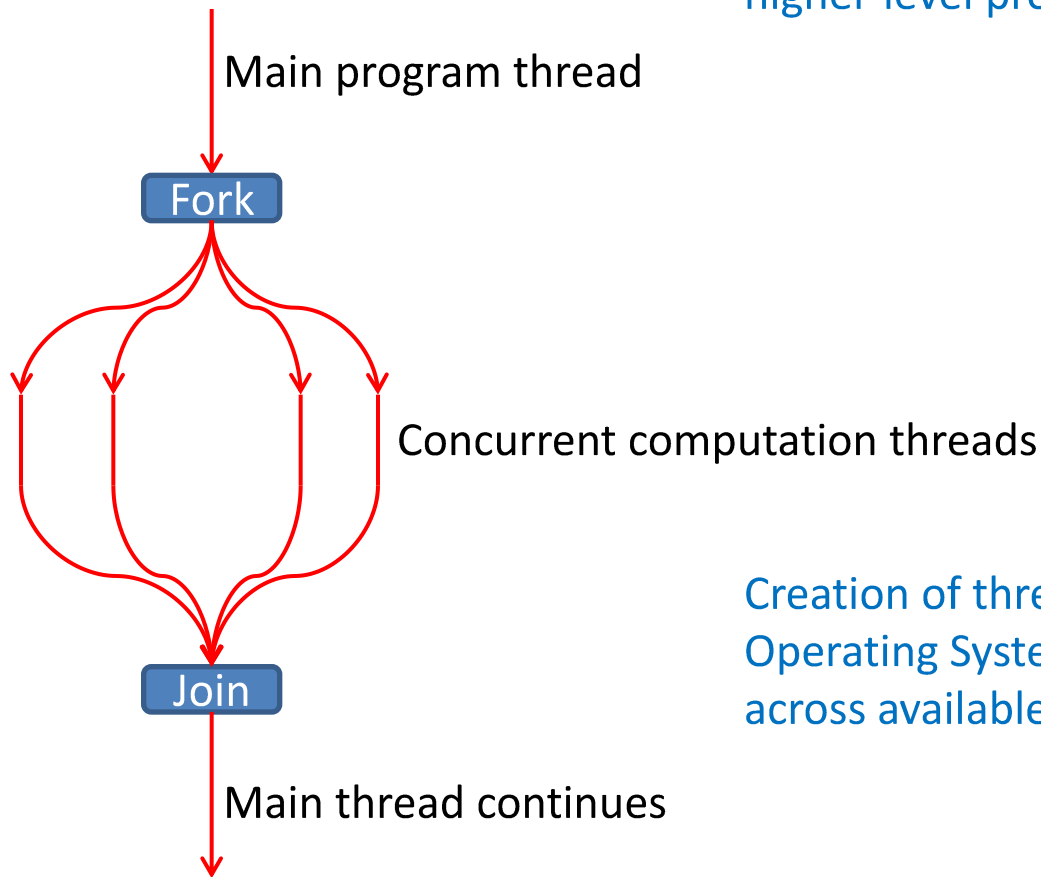
Presentation index

- Introduction
- Multi-threaded concurrency:
Data- versus Task-partitioning
- Parallelization with dependencies:
Reduction expressions or Streaming
- Multi-threading: difficult...
- Android: help from Pareon and Perf
- Conclusion



Creating multi-threaded concurrency

Basic fork-join pattern, created through different higher-level programming constructs



Creation of threads is application responsibility. Operating System handles run-time scheduling across available processors!

Parallelization – two partitioning options

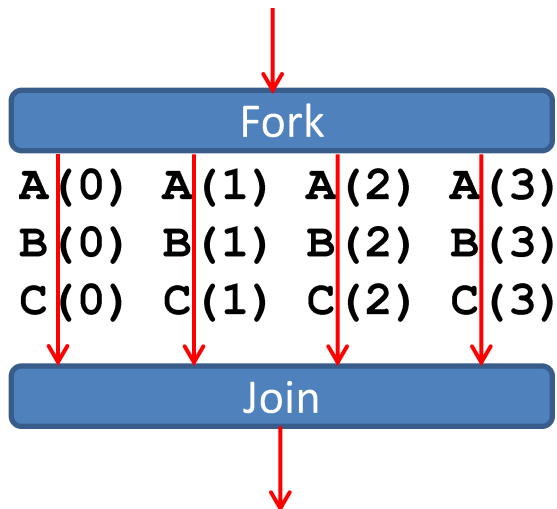
Source code:

```
for (i=0; i<4; i++) {  
    A(i);  
    B(i);  
    C(i);  
}
```

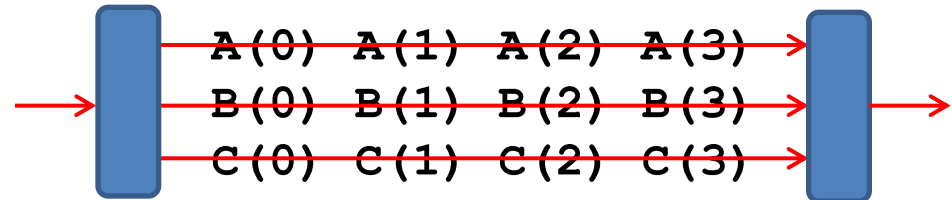
Sequential execution order:

A(0) A(1) A(2) A(3)
B(0) B(1) B(2) B(3)
C(0) C(1) C(2) C(3)

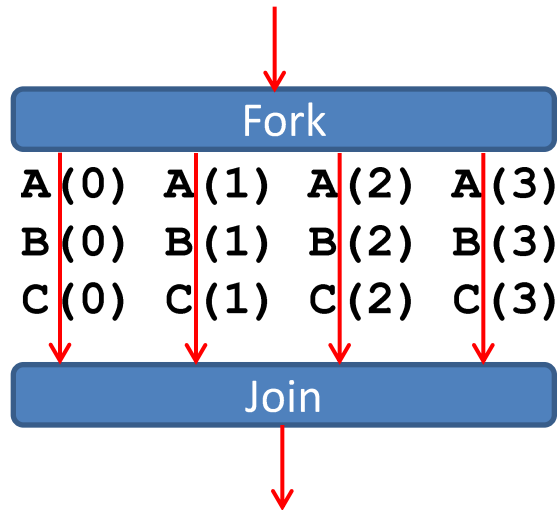
Data partitioning:



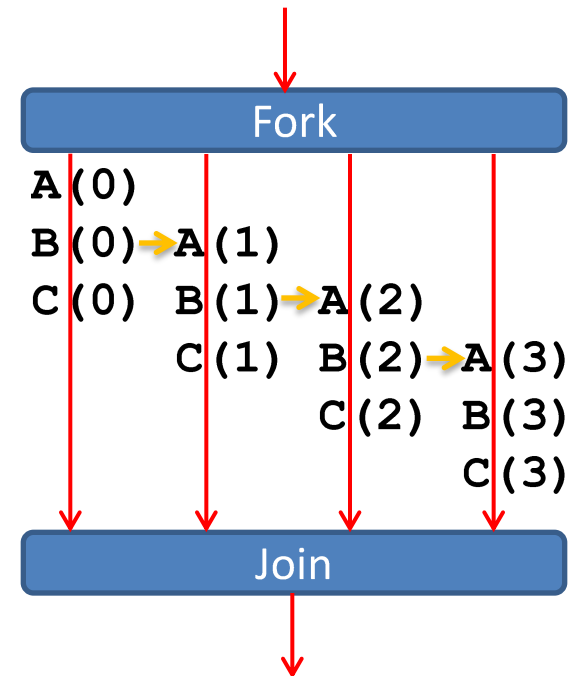
Task partitioning:



Issue: Data dependencies



Maybe, $B(i)$
produces a value
that is used by
 $A(i+1) \dots$



Adjust program source for parallelization:

- When feasible, remove inter-thread data dependencies
- Implement required data synchronization

Example Data dependencies

Variable assigned in loop body, used in later iteration

```
// search linked-list for matching items
// save matches in 'found' array of pointers
for (p = head, n_found = 0; p; p = p->next)
    if (match_criterion(p))
        found[n_found++] = p;
```

Cannot (easily/trivially) spawn data-parallel tasks!

- No direct parallel access to list members ***p**
- No direct way to assign index to matched item **n_found**
- Maybe more problems hidden in **match_criterion()**

Presentation index

- Introduction
- Multi-threaded concurrency:
Data- versus Task-partitioning
- Parallelization with dependencies:
Reduction expressions or Streaming
- Multi-threading: difficult...
- Android: help from Pareon and Perf
- Conclusion



Can do: reduction data dependencies

- Reduction expressions: accumulate results of loop bodies with commutative operations
- Freedom of re-ordering allows to break sequential constraints

```
// conditionally accumulate results
```

```
int acc = 0;  
for (i=0; i<N; i++)  
{  
    int result = some_work(i);  
    if (some_condition(i))  
        acc += result;  
}
```

```
...use of acc ...
```

- Commutative operations are basic math like +, *, &&, &, ||, but also more complex operations like *'add item to set'*.
- Three(?) different methods to handle these ...

Three methods for reduction dependencies

- Create thread-local copies of the accumulator. Accumulate over local copy in each thread. Merge the partial accumulators after thread-join. Eg. created automatically by:
`#pragma omp parallel for reduction(...)`
- Maintain single accumulator, synchronize updates through atomic operations. Eg. in C11 or C++11:
`atomic_add_fetch(&acc, result);`
`std::atomic<int> acc;`
`acc += result;`
- Maintain single accumulator, synchronize updates through protection by acquiring and releasing semaphores. Eg. Used by Intel “Threaded Building Blocks” (C++):
`concurrent_unordered_set<...> s;`
`s.insert(...);`

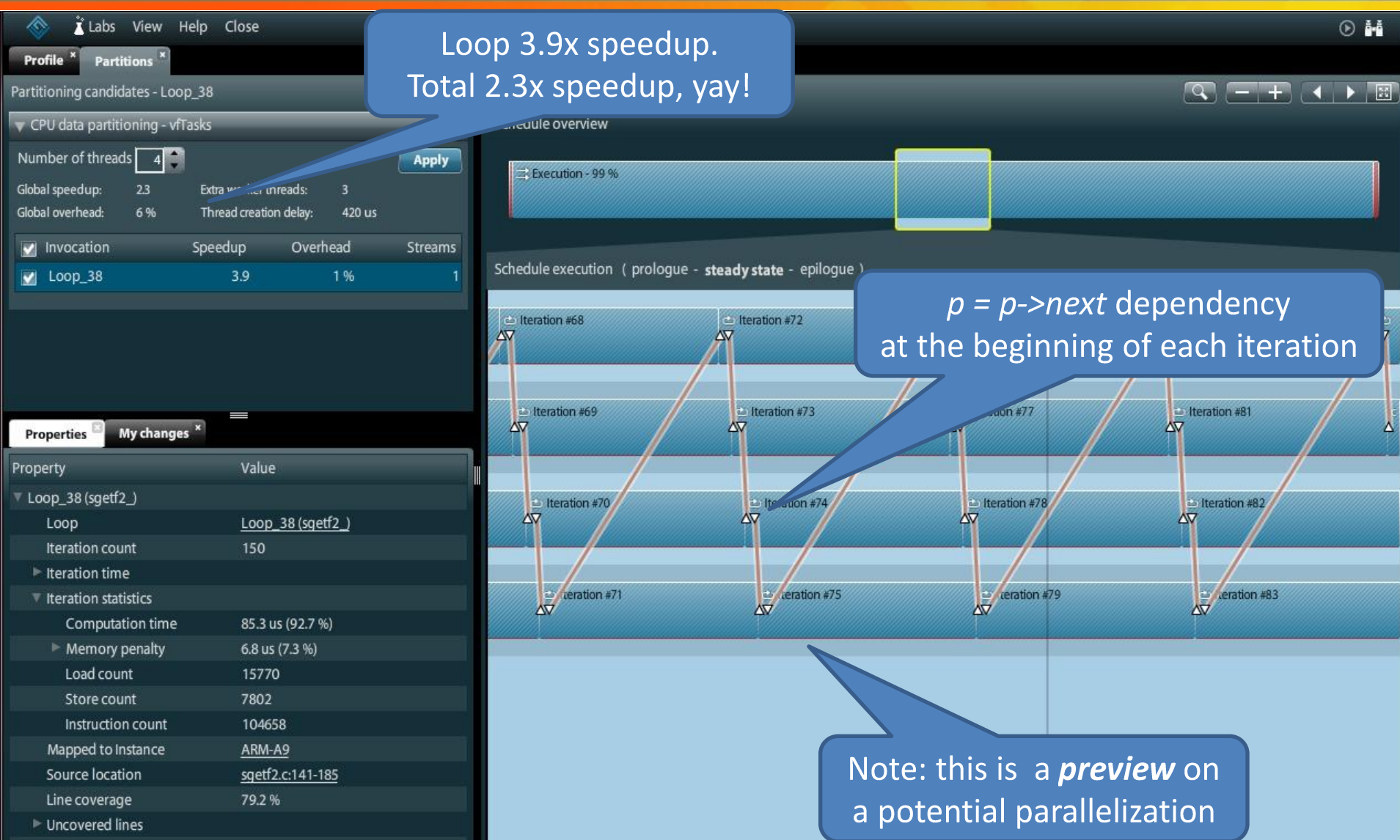
PAREON: Parallelization Analysis – 1

The screenshot shows the PAREON tool interface with several panels and callouts:

- Properties Panel (Top Left):** Displays memory dependency cluster 278 (carried by Loop_16429). It lists properties like cluster context, #dependencies, ignore with data, #loop-carried transfers, and memory dependency.
- Dynamic Help / Cheat sheets Panel (Top Right):** Provides an introduction to the tool, explaining the process of parallelizing a program and the use of help and check icons.
- Profile / Partitions Panel (Bottom Left):** A table showing the coverage and delay for various tasks and loops. The data is as follows:

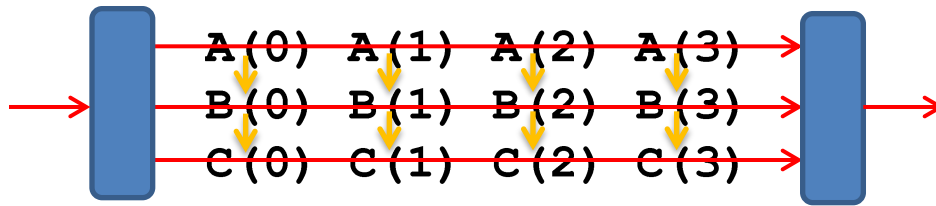
Name	Coverage	Delay
Parallel_16428	97	94.06
Task_16429	97	94.06
Mat::ptr	100	0.00
Parallel_16430	100	27.44
Task_16431	100	27.43
Loop_4178	100	0.34
_vf_memset	0	0.00
Loop_4179	100	48.09
Loop_4215	0	0.00
Loop_4216	0	0.00
Parallel_16432	0	0.00
Loop_4217	0	0.00
Loop_4218	0	0.00
Loop_4219	0	0.00
Loop_4220	0	0.00
Loop_4221	0	0.00
Loop_4222	0	0.00
Loop_4223	0	0.00
Loop_4224	0	0.00
Loop_4225	0	0.00
Loop_4226	0	0.00
Loop_4227	0	0.00
Loop_4228	0	0.00
Loop_4229	0	0.00
Loop_4230	0	0.00
Loop_4231	0	0.00
Loop_4232	0	0.00
Loop_4233	0	0.00
Loop_4234	0	0.00
Loop_4235	0	0.00
Loop_4236	0	0.00
Loop_4237	0	0.00
Loop_4238	0	0.00
Loop_4239	0	0.00
Loop_4240	0	0.00
Loop_4241	0	0.00
Loop_4242	0	0.00
Loop_4243	0	0.00
Loop_4244	0	0.00
Loop_4245	0	0.00
Loop_4246	0	0.00
Loop_4247	0	0.00
Loop_4248	0	0.00
Loop_4249	0	0.00
Loop_4250	0	0.00
Loop_4251	0	0.00
Loop_4252	0	0.00
Loop_4253	0	0.00
Loop_4254	0	0.00
Loop_4255	0	0.00
Loop_4256	0	0.00
Loop_4257	0	0.00
Loop_4258	0	0.00
Loop_4259	0	0.00
Loop_4260	0	0.00
Loop_4261	0	0.00
Loop_4262	0	0.00
Loop_4263	0	0.00
Loop_4264	0	0.00
Loop_4265	0	0.00
Loop_4266	0	0.00
Loop_4267	0	0.00
Loop_4268	0	0.00
Loop_4269	0	0.00
Loop_4270	0	0.00
Loop_4271	0	0.00
Loop_4272	0	0.00
Loop_4273	0	0.00
Loop_4274	0	0.00
Loop_4275	0	0.00
Loop_4276	0	0.00
Loop_4277	0	0.00
Loop_4278	0	0.00
Loop_4279	0	0.00
Loop_4280	0	0.00
Loop_4281	0	0.00
Loop_4282	0	0.00
Loop_4283	0	0.00
Loop_4284	0	0.00
Loop_4285	0	0.00
Loop_4286	0	0.00
Loop_4287	0	0.00
Loop_4288	0	0.00
Loop_4289	0	0.00
Loop_4290	0	0.00
Loop_4291	0	0.00
Loop_4292	0	0.00
Loop_4293	0	0.00
Loop_4294	0	0.00
Loop_4295	0	0.00
Loop_4296	0	0.00
Loop_4297	0	0.00
Loop_4298	0	0.00
Loop_4299	0	0.00
Loop_4300	0	0.00
Loop_4301	0	0.00
Loop_4302	0	0.00
Loop_4303	0	0.00
Loop_4304	0	0.00
Loop_4305	0	0.00
Loop_4306	0	0.00
Loop_4307	0	0.00
Loop_4308	0	0.00
Loop_4309	0	0.00
Loop_4310	0	0.00
Loop_4311	0	0.00
Loop_4312	0	0.00
Loop_4313	0	0.00
Loop_4314	0	0.00
Loop_4315	0	0.00
Loop_4316	0	0.00
Loop_4317	0	0.00
Loop_4318	0	0.00
Loop_4319	0	0.00
Loop_4320	0	0.00
Loop_4321	0	0.00
Loop_4322	0	0.00
Loop_4323	0	0.00
Loop_4324	0	0.00
Loop_4325	0	0.00
Loop_4326	0	0.00
Loop_4327	0	0.00
Loop_4328	0	0.00
Loop_4329	0	0.00
Loop_4330	0	0.00
Loop_4331	0	0.00
Loop_4332	0	0.00
Loop_4333	0	0.00
Loop_4334	0	0.00
Loop_4335	0	0.00
Loop_4336	0	0.00
Loop_4337	0	0.00
Loop_4338	0	0.00
Loop_4339	0	0.00
Loop_4340	0	0.00
Loop_4341	0	0.00
Loop_4342	0	0.00
Loop_4343	0	0.00
Loop_4344	0	0.00
Loop_4345	0	0.00
Loop_4346	0	0.00
Loop_4347	0	0.00
Loop_4348	0	0.00
Loop_4349	0	0.00
Loop_4350	0	0.00
Loop_4351	0	0.00
Loop_4352	0	0.00
Loop_4353	0	0.00
Loop_4354	0	0.00
Loop_4355	0	0.00
Loop_4356	0	0.00
Loop_4357	0	0.00
Loop_4358	0	0.00
Loop_4359	0	0.00
Loop_4360	0	0.00
Loop_4361	0	0.00
Loop_4362	0	0.00
Loop_4363	0	0.00
Loop_4364	0	0.00
Loop_4365	0	0.00
Loop_4366	0	0.00
Loop_4367	0	0.00
Loop_4368	0	0.00
Loop_4369	0	0.00
Loop_4370	0	0.00
Loop_4371	0	0.00
Loop_4372	0	0.00
Loop_4373	0	0.00
Loop_4374	0	0.00
Loop_4375	0	0.00
Loop_4376	0	0.00
Loop_4377	0	0.00
Loop_4378	0	0.00
Loop_4379	0	0.00
Loop_4380	0	0.00
Loop_4381	0	0.00
Loop_4382	0	0.00
Loop_4383	0	0.00
Loop_4384	0	0.00
Loop_4385	0	0.00
Loop_4386	0	0.00
Loop_4387	0	0.00
Loop_4388	0	0.00
Loop_4389	0	0.00
Loop_4390	0	0.00
Loop_4391	0	0.00
Loop_4392	0	0.00
Loop_4393	0	0.00
Loop_4394	0	0.00
Loop_4395	0	0.00
Loop_4396	0	0.00
Loop_4397	0	0.00
Loop_4398	0	0.00
Loop_4399	0	0.00
Loop_4400	0	0.00
Loop_4401	0	0.00
Loop_4402	0	0.00
Loop_4403	0	0.00
Loop_4404	0	0.00
Loop_4405	0	0.00
Loop_4406	0	0.00
Loop_4407	0	0.00
Loop_4408	0	0.00
Loop_4409	0	0.00
Loop_4410	0	0.00
Loop_4411	0	0.00
Loop_4412	0	0.00
Loop_4413	0	0.00
Loop_4414	0	0.00
Loop_4415	0	0.00
Loop_4416	0	0.00
Loop_4417	0	0.00
Loop_4418	0	0.00
Loop_4419	0	0.00
Loop_4420	0	0.00
Loop_4421	0	0.00
Loop_4422	0	0.00
Loop_4423	0	0.00
Loop_4424	0	0.00
Loop_4425	0	0.00
Loop_4426	0	0.00
Loop_4427	0	0.00
Loop_4428	0	0.00
Loop_4429	0	0.00
Loop_4430	0	0.00
Loop_4431	0	0.00
Loop_4432	0	0.00
Loop_4433	0	0.00
Loop_4434	0	0.00
Loop_4435	0	0.00
Loop_4436	0	0.00
Loop_4437	0	0.00
Loop_4438	0	0.00
Loop_4439	0	0.00
Loop_4440	0	0.00
Loop_4441	0	0.00
Loop_4442	0	0.00
Loop_4443	0	0.00
Loop_4444	0	0.00
Loop_4445	0	0.00
Loop_4446	0	0.00
Loop_4447	0	0.00
Loop_4448	0	0.00
Loop_4449	0	0.00
Loop_4450	0	0.00
Loop_4451	0	0.00
Loop_4452	0	0.00
Loop_4453	0	0.00
Loop_4454	0	0.00
Loop_4455	0	0.00
Loop_4456	0	0.00
Loop_4457	0	0.00
Loop_4458	0	0.00
Loop_4459	0	0.00
Loop_4460	0	0.00
Loop_4461	0	0.00
Loop_4462	0	0.00
Loop_4463	0	0.00
Loop_4464	0	0.00
Loop_4465	0	0.00
Loop_4466	0	0.00
Loop_4467	0	0.00
Loop_4468	0	0.00
Loop_4469	0	0.00
Loop_4470	0	0.00
Loop_4471	0	0.00
Loop_4472	0	0.00
Loop_4473	0	0.00
Loop_4474	0	0.00
Loop_4475	0	0.00
Loop_4476	0	0.00
Loop_4477	0	0.00
Loop_4478	0	0.00
Loop_4479	0	0.00
Loop_4480	0	0.00
Loop_4481	0	0.00
Loop_4482	0	0.00
Loop_4483	0	0.00
Loop_4484	0	0.00
Loop_4485	0	0.00
Loop_4486	0	0.00
Loop_4487	0	0.00
Loop_4488	0	0.00
Loop_4489	0	0.00
Loop_4490	0	0.00
Loop_4491	0	0.00
Loop_4492	0	0.00
Loop_4493	0	0.00
Loop_4494	0	0.00
Loop_4495	0	0.00
Loop_4496	0	0.00
Loop_4497	0	0.00
Loop_4498	0	0.00
Loop_4499	0	0.00
Loop_4500	0	0.00
Loop_4501	0	0.00
Loop_4502	0	0.00
Loop_4503	0	0.00
Loop_4504	0	0.00
Loop_4505	0	0.00
Loop_4506	0	0.00
Loop_4507	0	0.00
Loop_4508	0	0.00
Loop_4509	0	0.00
Loop_4510	0	0.00
Loop_4511	0	0.00
Loop_4512	0	0.00
Loop_4513	0	0.00
Loop_4514	0	0.00
Loop_4515	0	0.00
Loop_4516	0	0.00
Loop_4517	0	0.00
Loop_4518	0	0.00
Loop_4519	0	0.00
Loop_4520	0	0.00
Loop_4521	0	0.00
Loop_4522	0	0.00
Loop_4523	0	0.00
Loop_4524	0	0.00
Loop_4525	0	0.00
Loop_4526	0	0.00
Loop_4527	0	0.00
Loop_4528	0	0.00
Loop_4529	0	0.00
Loop_4530	0	0.00
Loop_4531	0	0.00
Loop_4532	0	0.00
Loop_4533	0	0.00
Loop_4534	0	0.00
Loop_4535	0	0.00
Loop_4536	0	0.00
Loop_4537	0	0.00
Loop_4538	0	0.00
Loop_4539	0	0.00
Loop_4540	0	0.00
Loop_4541	0	0.00
Loop_4542	0	0.00
Loop_4543	0	0.00
Loop_4544	0	0.00
Loop_4545	0	0.00
Loop_4546	0	0.00
Loop_4547	0	0.00
Loop_4548	0	0.00
Loop_4549	0	0.00
Loop_4550	0	0.00
Loop_4551	0	0.00
Loop_4552	0	0.00
Loop_4553	0	0.00
Loop_4554	0	0.00
Loop_4555	0	0.00
Loop_4556	0	0.00
Loop_4557	0	0.00
Loop_4558	0	0.00
Loop_4559	0	0.00
Loop_4560	0	0.00
Loop_4561	0	0.00
Loop_4562	0	0.00
Loop_4563	0	0.00
Loop_4564	0	0.00
Loop_4565	0	0.00
Loop_4566	0	0.00
Loop_4567	0	0.00
Loop_4568	0	0.00
Loop_4569	0	0.00
Loop_4570	0	0.00
Loop_4571	0	0.00
Loop_4572	0	0.00
Loop_4573	0	0.00
Loop_4574	0	0.00
Loop_4575	0	0.00
Loop_4576	0	0.00
Loop_4577	0	0.00
Loop_4578	0	0.00
Loop_4579	0	0.00
Loop_4580	0	0.00
Loop_4581	0	0.00
Loop_4582	0	0.00
Loop_4583	0	0.00
Loop_4584	0	0.00
Loop_4585	0	0.00
Loop_4586	0	0.00
Loop_4587	0	0.00
Loop_4588	0	0.00
Loop_4589	0	0.00
Loop_4590	0	0.00
Loop_4591	0	0.00
Loop_4592	0	0.00
Loop_4593	0	0.00
Loop_4594	0	0.00
Loop_4595	0	0.00
Loop_4596	0	0.00
Loop_4597	0	0.00
Loop_4598	0	0.00
Loop_4599	0	0.00
Loop_4600	0	0.00
Loop_4601	0	0.00
Loop_4602	0	0.00
Loop_4603	0	0.00
Loop_4604	0	0.00
Loop_4605	0	0.00
Loop_4606	0	0.00
Loop_4607	0	0.00
Loop_4608	0	0.

PAREON: Parallelization Analysis - 2

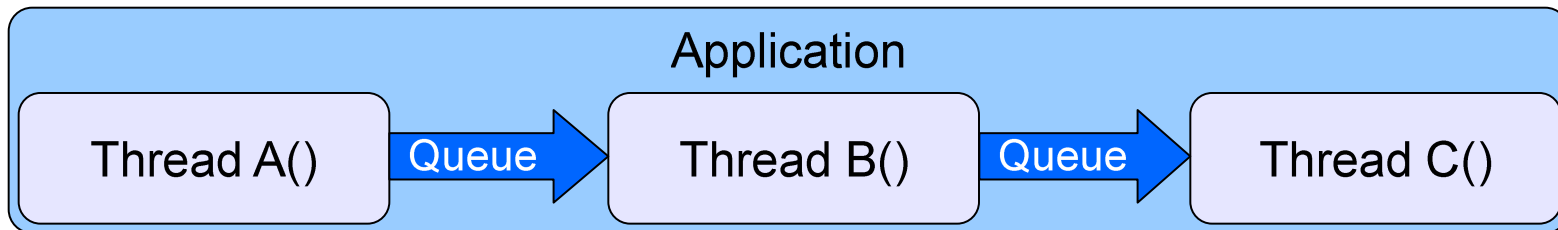


Pipelining: Data deps & functional partitioning

Functional partitioning with inter-thread dependencies:



Producer-Consumer pattern:



Queue implementation solves dependencies:

- **Solve Data dependencies:** Consumer thread waits for available data (stalls until queue is non-empty)
- **Solve Anti dependencies:** Producer thread creates next item in next memory location (prevents overwriting previous value)

Presentation index

- Introduction
- Multi-threaded concurrency:
Data- versus Task-partitioning
- Parallelization with dependencies:
Reduction expressions or Streaming
- **Multi-threading: difficult...**
- **Android: help from Pareon and Perf**
- Conclusion



Concurrent C/C++ programming: Pitfalls

Risc introduction of functional errors:

- Overlooking use of shared/global variables
(deep down inside called functions, or inside 3rd party library)
- Overlooking exceptions that are raised and caught outside studied scope
- Incorrect use of semaphores: flawed protection, deadlocks

Unexpected performance issues:

- Underestimation of time spent in added multi-threading or synchronization code and libraries
- Underestimation of other penalties in OS and HW
(inter-core cache penalties, context switches, clock-frequency reductions)

Parallel programming remains hard!

Development of parallel code

Guidelines:

- Base upon a sequential program:
functional and performance reference
- Apply higher-level parallelization patterns and primitives:
clear semantics, re-use code, reduce risk
- Use tooling for analysis and verification
 - Prevent introduction of hard-to-find bugs
 - Prevent recoding effort that does not perform

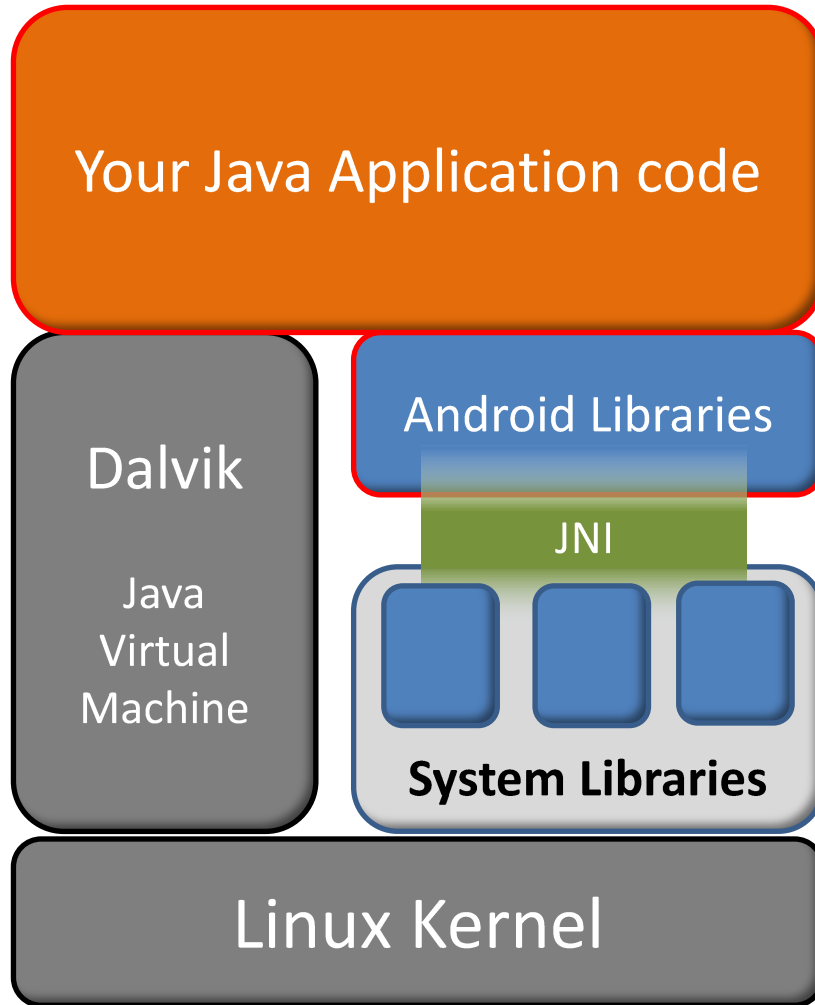
Managable development process!

Presentation index

- Introduction
- Multi-threaded concurrency:
Data- versus Task-partitioning
- Parallelization with dependencies:
Reduction expressions or Streaming
- Multi-threading: difficult...
- **Android: help from Pareon and Perf**
- **Conclusion**



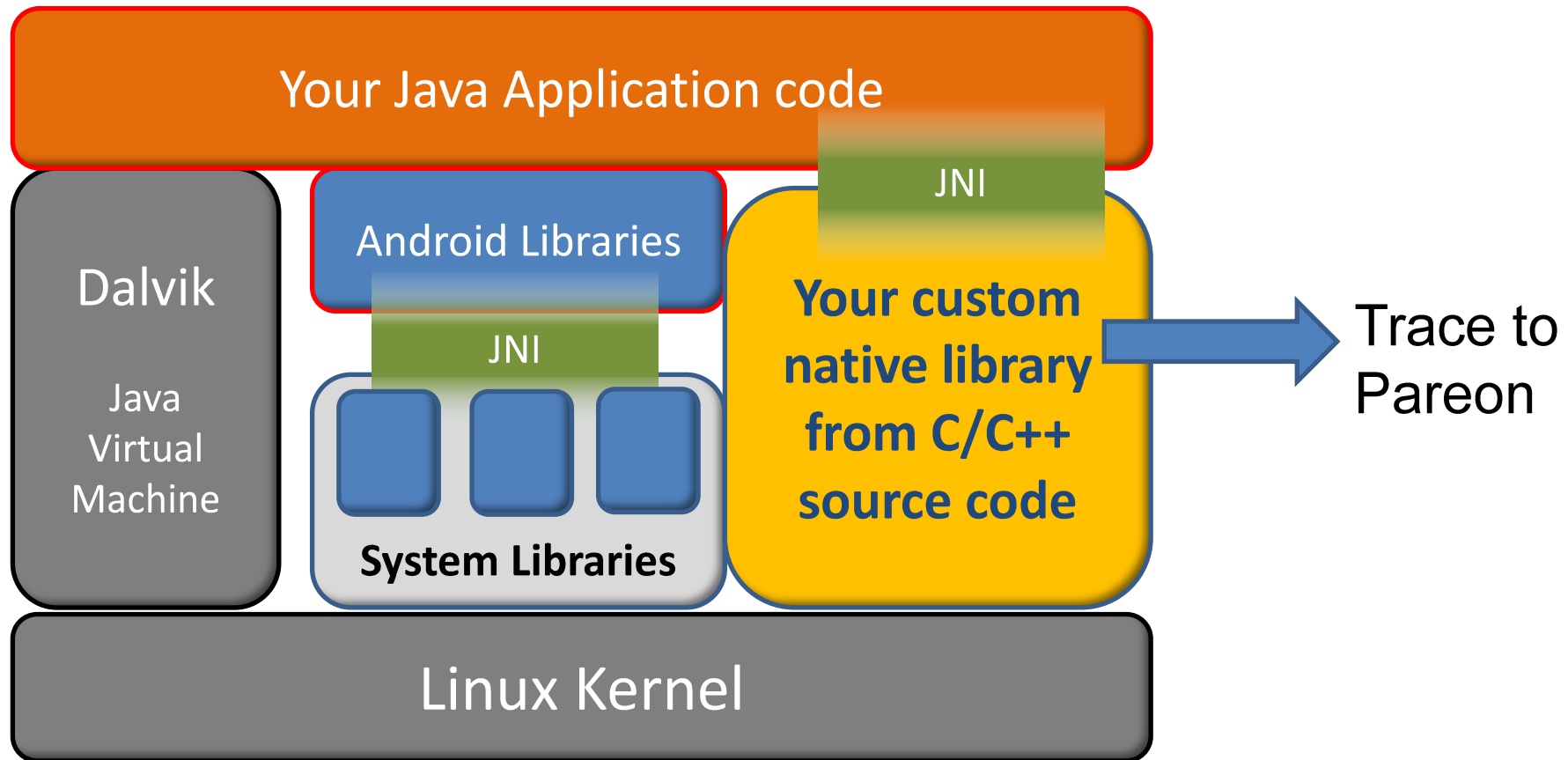
Android Application 1: Plain, just Java



Many apps have no critical CPU load
For now, no Java support in Pareon

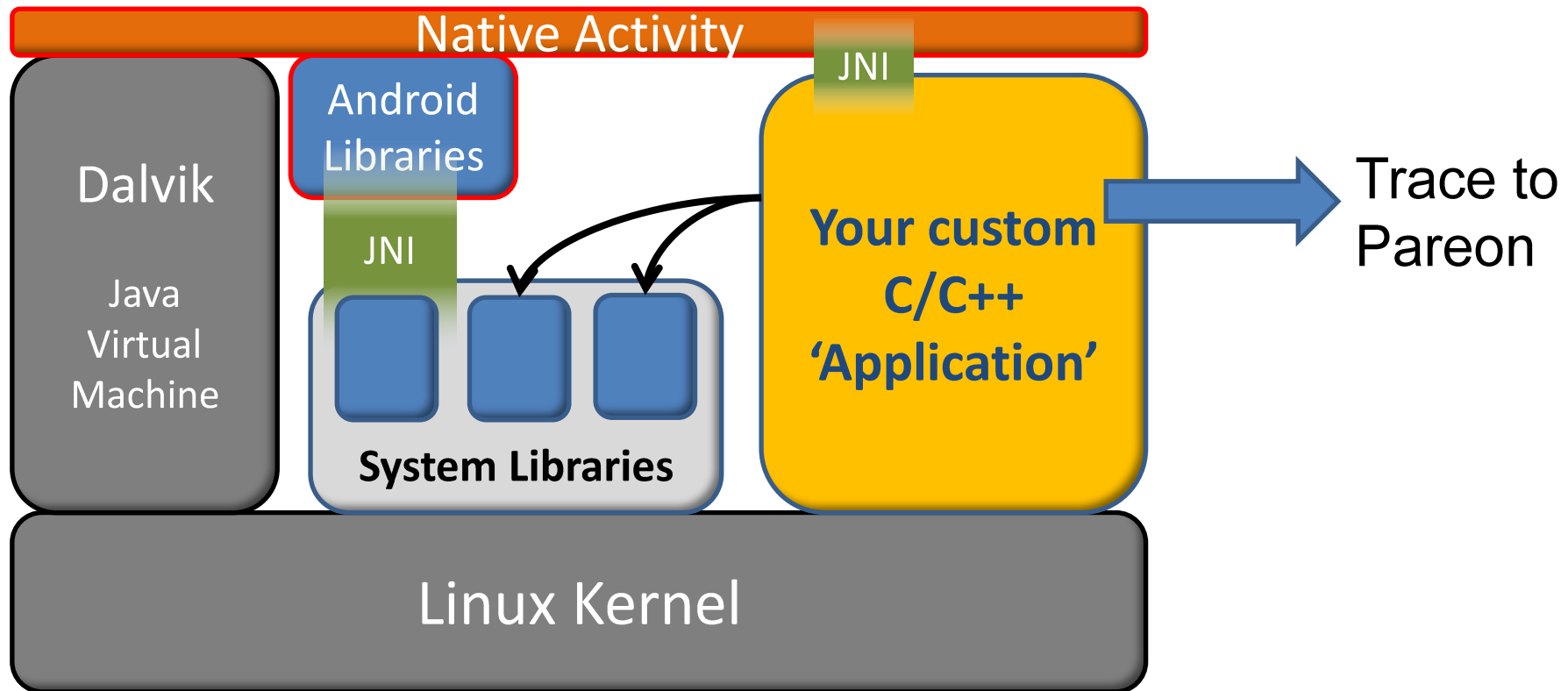
Android Application 2: with native libraries

Apps can include “native” binary code for best performance

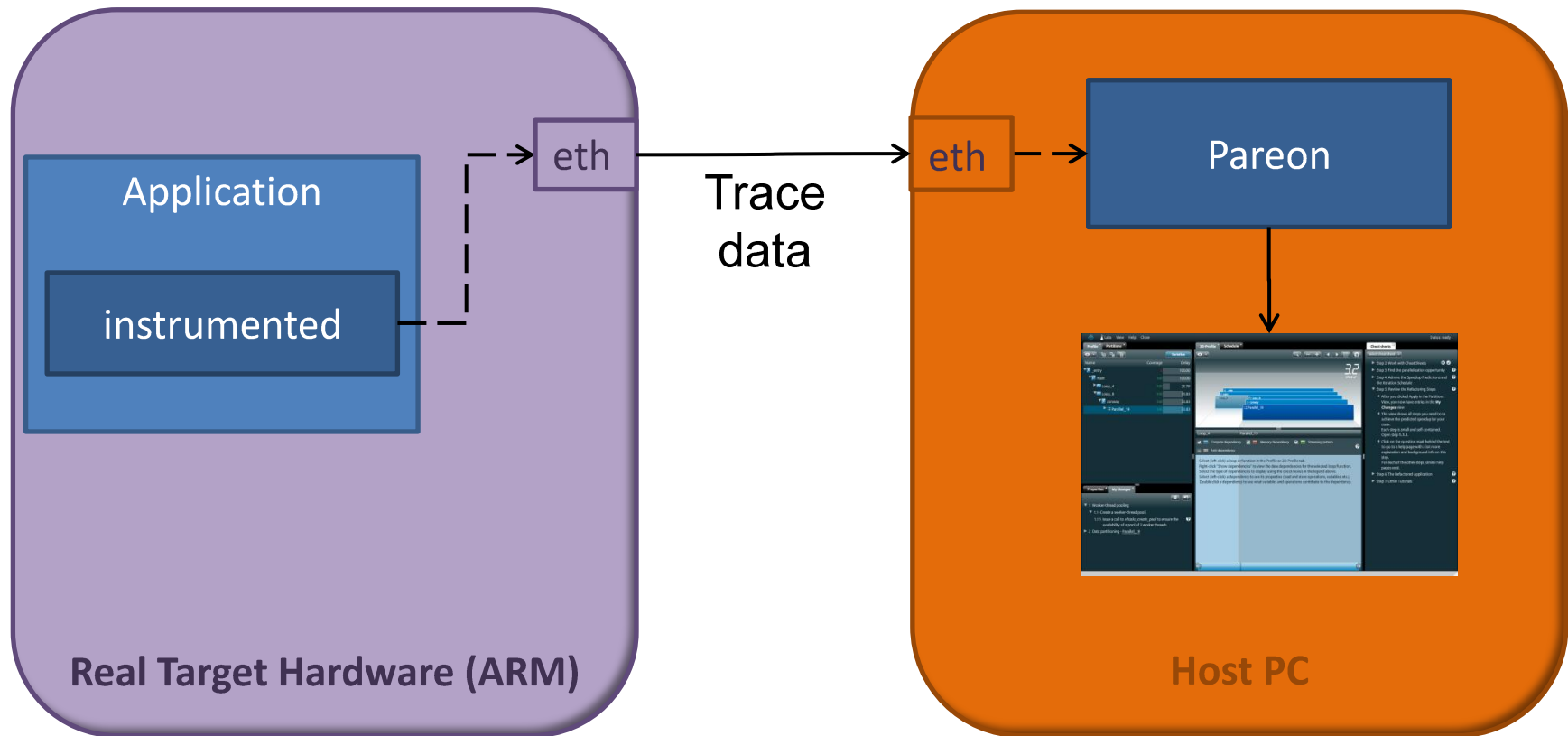


Android Application 3: NativeActivity

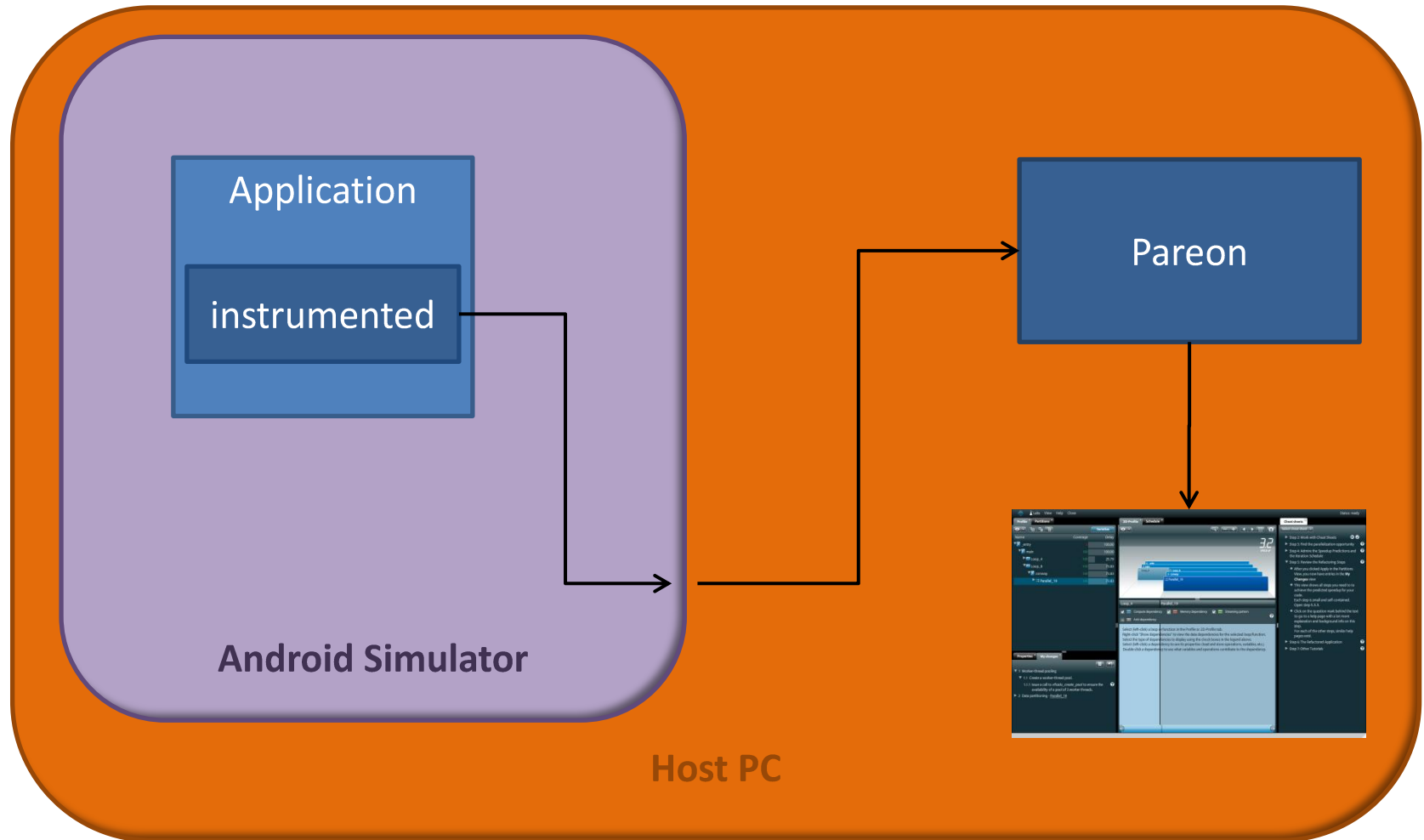
“Native activities” are created without Java source code



Application Analysis on Android target



System Setup using Android Simulator



NDK plasma demo app analyzed on Android



Finding data parallelism on Android

The screenshot displays the VectorFabrics Labs interface, which is used for analyzing and optimizing code for data parallelism. The interface is divided into several panels:

- Profile / Partitions:** This panel shows partitioning candidates for Loop_313. It includes a table of CPU data partitioning results.
- Properties / My changes:** This panel provides detailed information about the selected partitioning candidate, Loop_313.
- 2D-Profile / Schedule / plasma.c:** This panel shows a 2D profile of the code, highlighting the execution of Loop_313. It includes a speedup indicator and a visual representation of the code's execution flow.

Partitioning candidates - Loop_313

▼ CPU data partitioning - vfTasks

Number of threads: **Apply**

Global speedup: 2.4 Extra worker threads: 3
Global overhead: 1 % Thread creation delay: 420 us

<input checked="" type="checkbox"/>	Invocation	Speedup	Overhead	Streams
<input checked="" type="checkbox"/>	Loop_120	4.0	1 %	0
<input checked="" type="checkbox"/>	Loop_189	4.0	1 %	0
<input checked="" type="checkbox"/>	Loop_313	3.9	1 %	0

Properties / My changes

Property	Value
▼ Compute dependency 313.46	
Ignore with data partitioning	induction expression
▼ Source	
Operation (+)	Loop_313 (fill_plasma)
Location	plasma.c:211
► Destinations	
#loop-carried transfers	68.2 Ki transfers/s
Loop carried	yes

2D-Profile / Schedule / plasma.c

Speedup: 1.0

Loop_313 total loop carried transfer rate: 136 Ki transfers/s
0 streaming pattern clusters (0.0 transfers/s); 0 data dependency clusters (0.0 transfers/s);
2 compute dependencies (136 Ki transfers/s); 0 anti- and output dependency clusters

Finding data parallelism on Android

The screenshot displays the VectorFabrics Labs interface, which is used for analyzing and optimizing code for data parallelism. The interface is divided into several panels:

- Profile / Partitions:** This panel shows partitioning candidates for 'Loop_313'. It includes a table of candidates and a table of properties.
- 2D-Profile / Schedule:** This panel shows a 2D profile of the code, with a red circle highlighting the 'Loop_313' region.
- Properties / My changes:** This panel shows the properties of the selected partition, including its compute dependency, source, and destinations.

Partitioning candidates - Loop_313

Invocation	Speedup	Overhead	Streams
<input checked="" type="checkbox"/> Loop_120	4.0	1 %	0
<input checked="" type="checkbox"/> Loop_189	4.0	1 %	0
<input checked="" type="checkbox"/> Loop_313	3.9	1 %	0

Properties - My changes

Property	Value
Compute dependency	313.46
Ignore with data partitioning	induction expression
Source	
Operation (+)	Loop_313 (fill_plasma)
Location	plasma.c:211
Destinations	
#loop-carried transfers	68.2 Ki transfers/s
Loop carried	yes

2D-Profile / Schedule: plasma.c

The 2D profile shows a series of blue bars representing different code regions. A red circle highlights the 'Loop_313' region. The 'Schedule' panel shows the execution of 'Loop_313' with a speedup of 1.0. The 'Properties' panel shows the details of 'Loop_313', including its compute dependency, source, and destinations.

Loop_313 total loop carried transfer rate: 136 Ki transfers/s

0 streaming pattern clusters (0.0 transfers/s); 0 data dependency clusters (0.0 transfers/s);
2 compute dependencies (136 Ki transfers/s); 0 anti- and output dependency clusters

Finding data parallelism on Android

The screenshot displays the VectorFabrics Labs interface, which is used for analyzing and optimizing code for data parallelism. The interface is divided into several panels:

- Profile / Partitions:** This panel shows partitioning candidates for `Loop_313`. It includes a table of CPU data partitioning results and a table of properties for the selected partition.
- 2D-Profile / Schedule:** This panel shows a 2D profile of the code, with a focus on `Loop_313`. It includes a legend for dependencies and a detailed view of the loop's transfer rate and dependencies.

Partitioning candidates - Loop_313

▼ CPU data partitioning - vfTasks

Number of threads: **4** (circled in red) **Apply**

Global speedup: 2.4 Extra worker threads: 3
Global overhead: 1 % Thread creation delay: 420 us

✓	Invocation	Speedup	Overhead	Streams
✓	Loop_120	4.0	1 %	0
✓	Loop_189	4.0	1 %	0
✓	Loop_313	3.9	1 %	0

Properties / My changes

Property	Value
▼ Compute dependency 313.46	
Ignore with data partitioning	induction expression
▼ Source	
Operation (+)	Loop_313 (fill_plasma)
Location	plasma.c:211
► Destinations	
#loop-carried transfers	68.2 Ki transfers/s
Loop carried	yes

2D-Profile / Schedule: plasma.c

1.0 SPEED-UP

Legend: ☒ Compute dependency, ☒ Memory dependency, ☒ Streaming pattern, ☒ Anti-dependency

Loop_313 total loop carried transfer rate: 136 Ki transfers/s
0 streaming pattern clusters (0.0 transfers/s); 0 data dependency clusters (0.0 transfers/s);
2 compute dependencies (136 Ki transfers/s); 0 anti- and output dependency clusters

Not parallelized: JNI call to render frame

Profile **Partitions** **2D-Profile** **Schedule**

Profile

Name	Cover...	Delay
ANativeWindow_lock		6.05
__android_log_print		0.00
stats_startFrame	100	0.07
clock_gettime		0.23
fill_plasma	70	42.10
Parallel_370	64	42.10
ANativeWindow_unlockAndPost		26.17
stats_endFrame	93	1.61

Properties **My changes**

Property	Value
ANativeWindow_unlockAndPost	
Intrinsic	ANativeWindow_unlockAndPost
Invocation time	
Estimated	1.9 ms
Constraint	<delay> ns
Invocation statistics	
Mapped to Instance	ARM-A9
Call location	plasma.c:407

2D-Profile **Schedule**

2.4
SPEED-UP

ANati... P. en... AN... Parallel_370 ANativeWindow_unloc...

☒ Compute dependency ☒ Memory dependency ☒ Streaming pattern

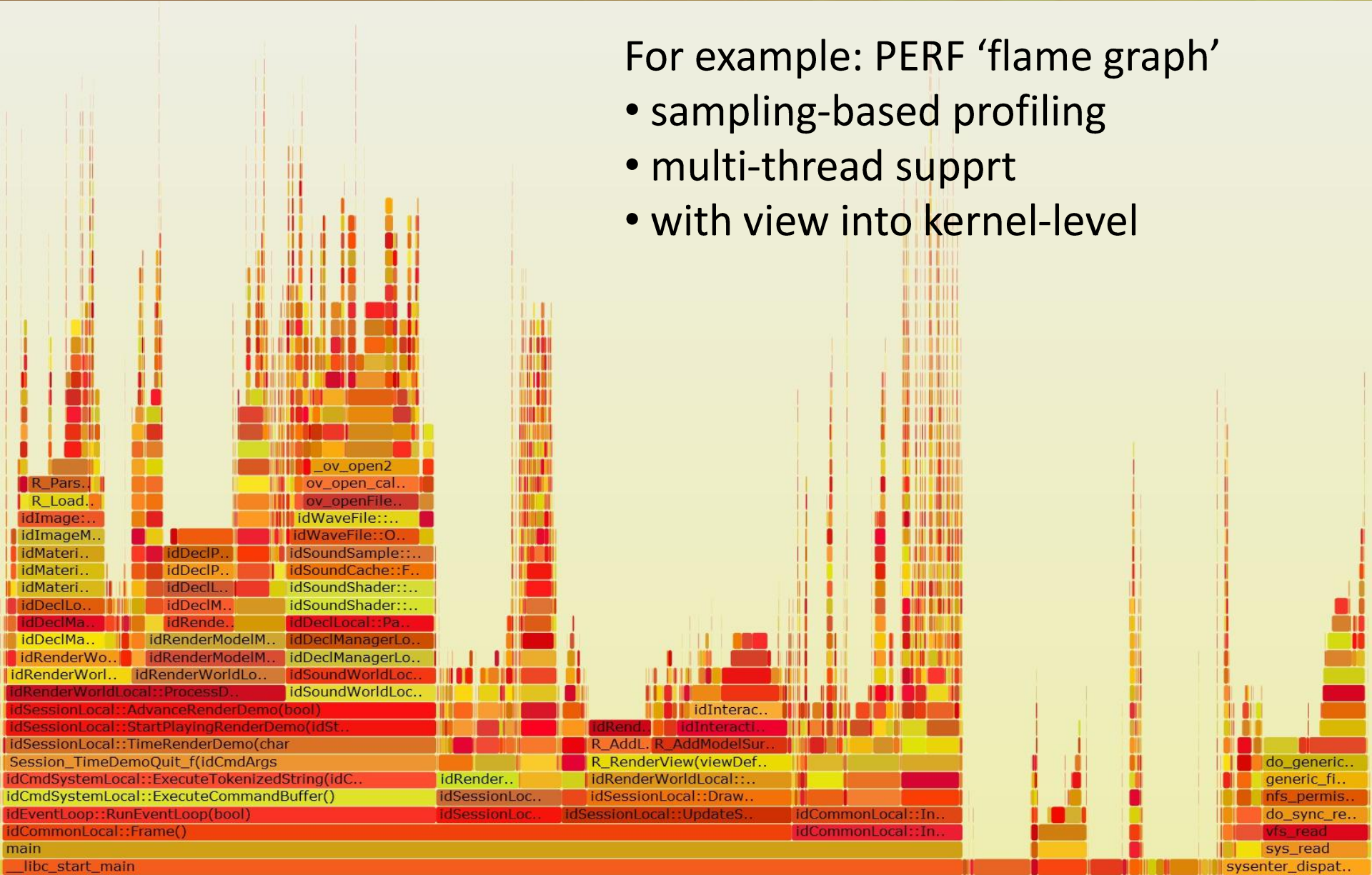
☒ Anti-dependency

Select (left-click) a loop or function in the Profile or 2D-Profile tab.
Right-click i.e. "Show internal dependencies" to view the internal data dependencies for the selected loop/function.
Select the type of dependencies to display using the check boxes in the legend above.
Select (left-click) a dependency to see its properties (load and store operations, variables, etc.).

Performance Verification

For example: PERF 'flame graph'

- sampling-based profiling
- multi-thread support
- with view into kernel-level



Conclusion

Today's gap:

- Multi-core CPUs are everywhere,
- Yet multi-threaded programming remains hard:
 - Risk of creating hard-to-locate bugs regarding dynamic data races and semaphore issues
 - Obtained speedup is lower than expected
- A sequential functional reference implementation ...
... helps to set a baseline for parallelization
- Android sets a new record in development complexity
- Proper tooling is needed to save on edit-verify development cycles

Today's gap:

- Multi-core CPUs are everywhere
- Yet multi-threaded programming remains hard:
 - Risk of creating hard-to-locate bugs regarding dynamic data races and semaphore issues
 - Obtained speedup is lower than expected
- A sequential functional reference implementation ...
... helps to set a baseline for parallelization
- Android sets a new record in development complexity
- Proper tooling is needed to save on edit-verify development cycles



Check www.vectorfabrics.com for a free demo on concurrency analysis



VectorFabrics

Thank you!