

FAASM: Lightweight Isolation for Efficient Stateful Serverless Computing

Simon Shillaker and Peter Pietzuch

Large-scale Data and Systems Group, Imperial College London

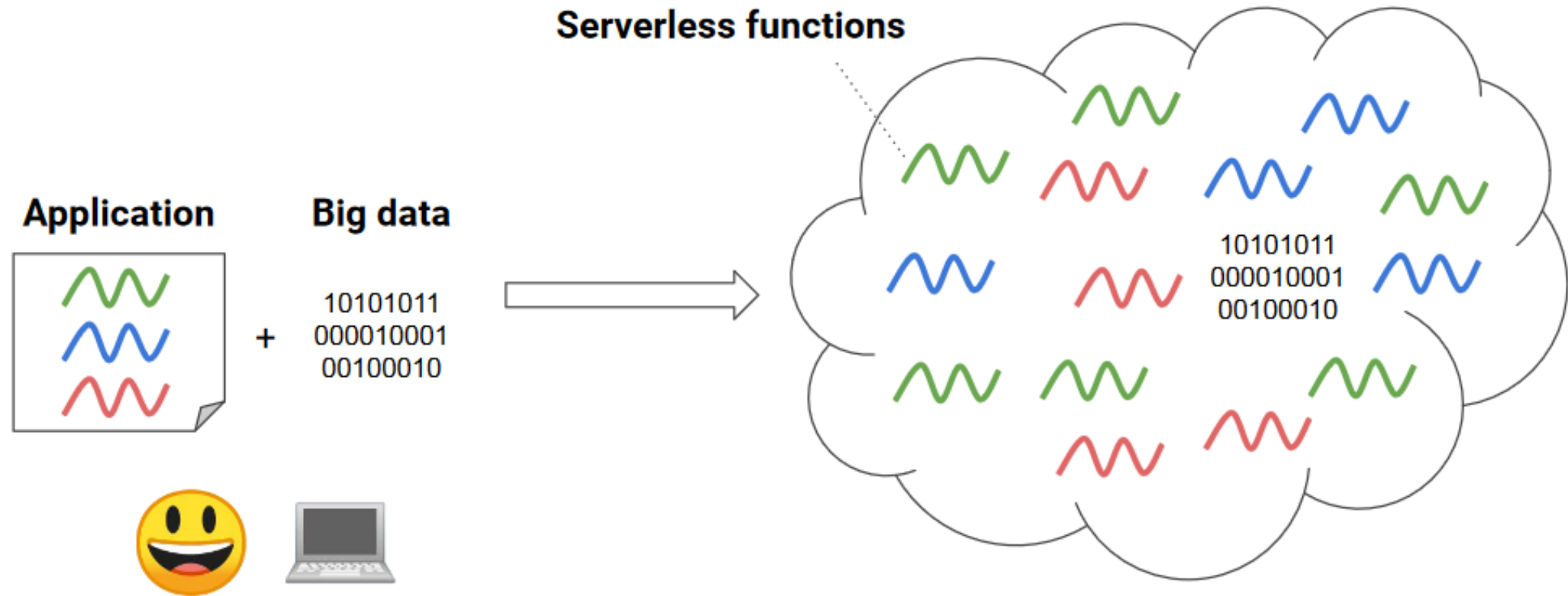


LSDS

Large-Scale Data & Systems Group

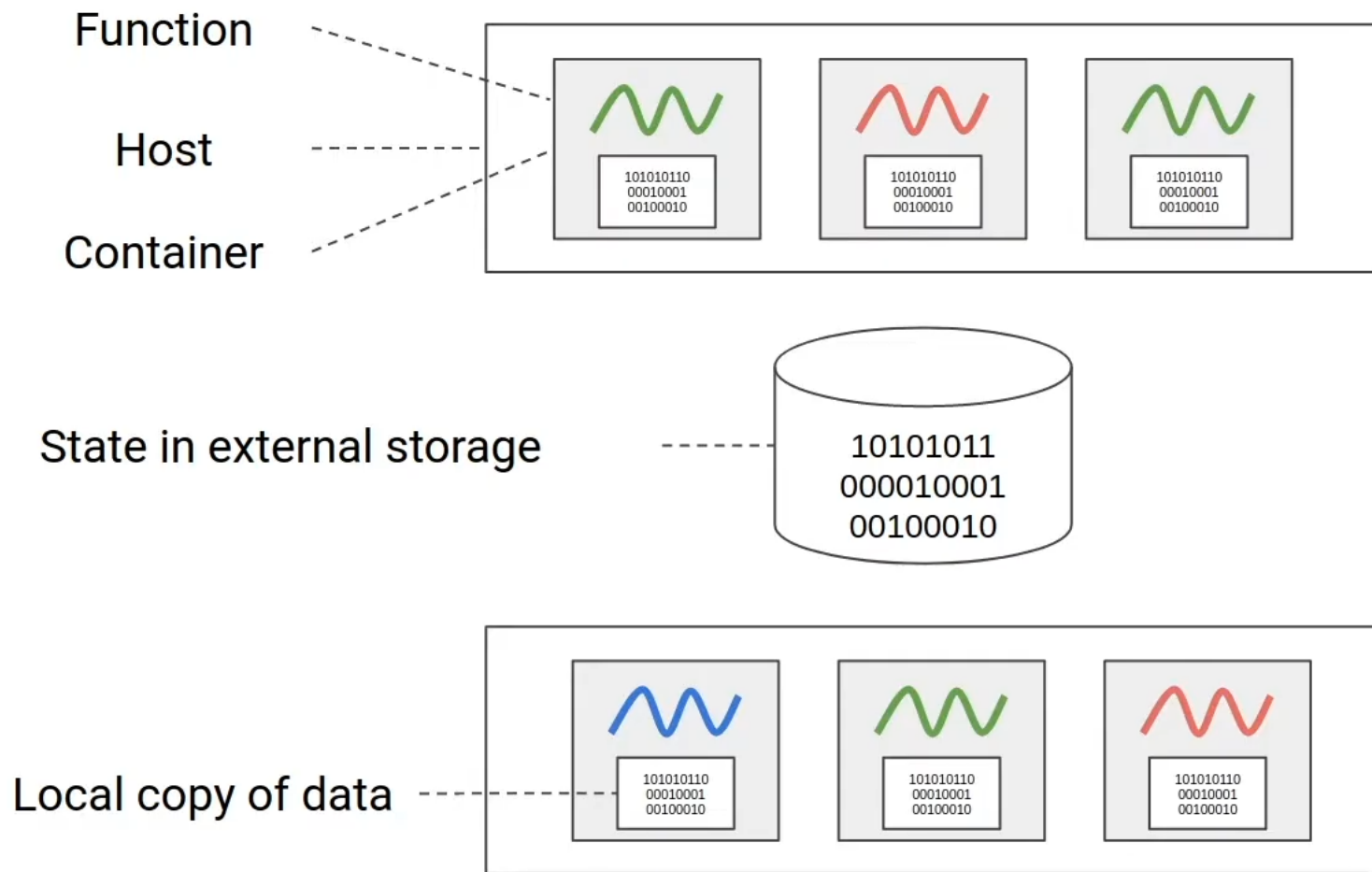
**Imperial College
London**

Serverless Big Data Vision

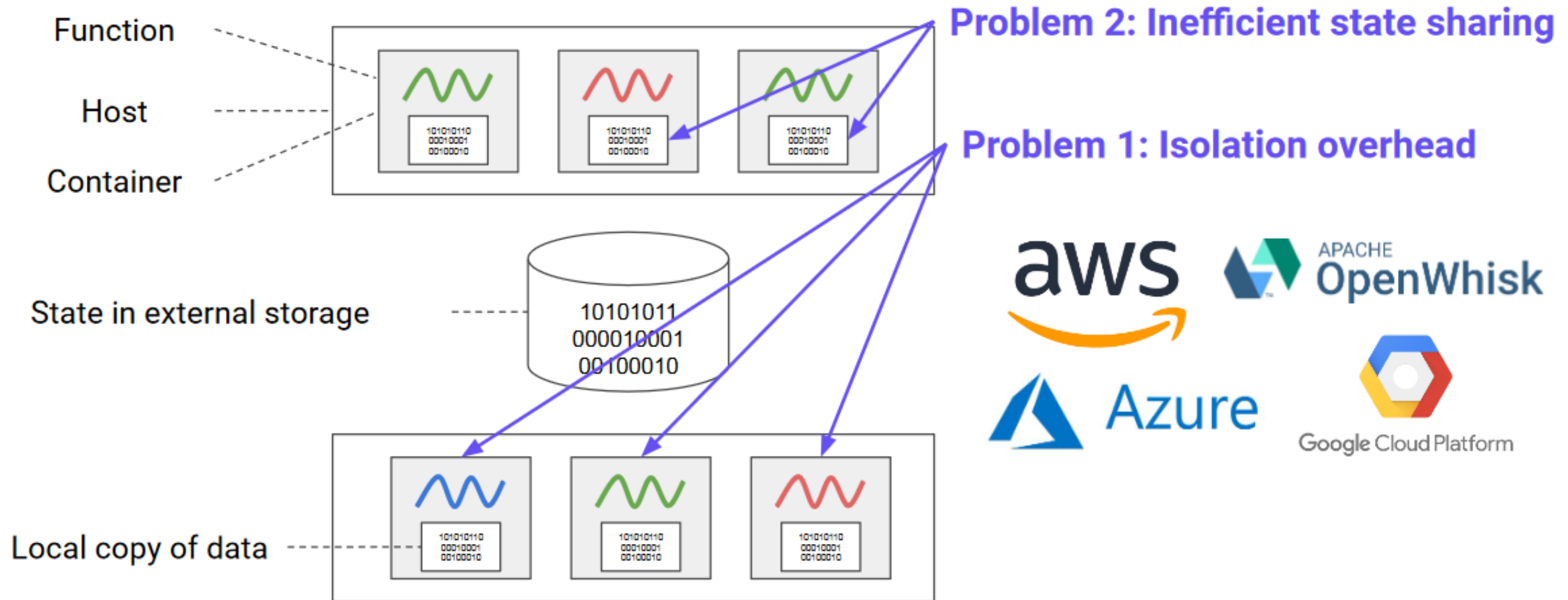


Cheap, highly scalable big data processing

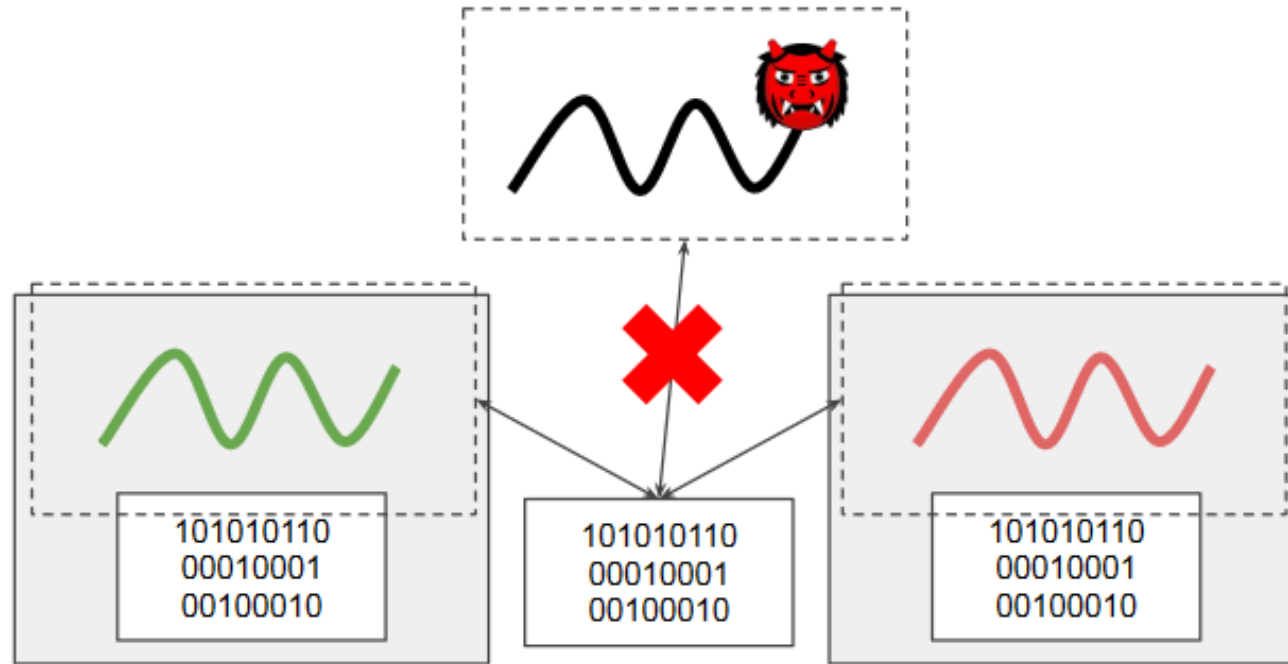
Serverless Under the Hood



Serverless Under the Hood



How Do We Efficiently Share State But Maintain Isolation?



**We need an isolation mechanism
that gives us fine-grained control over
memory**

WebAssembly

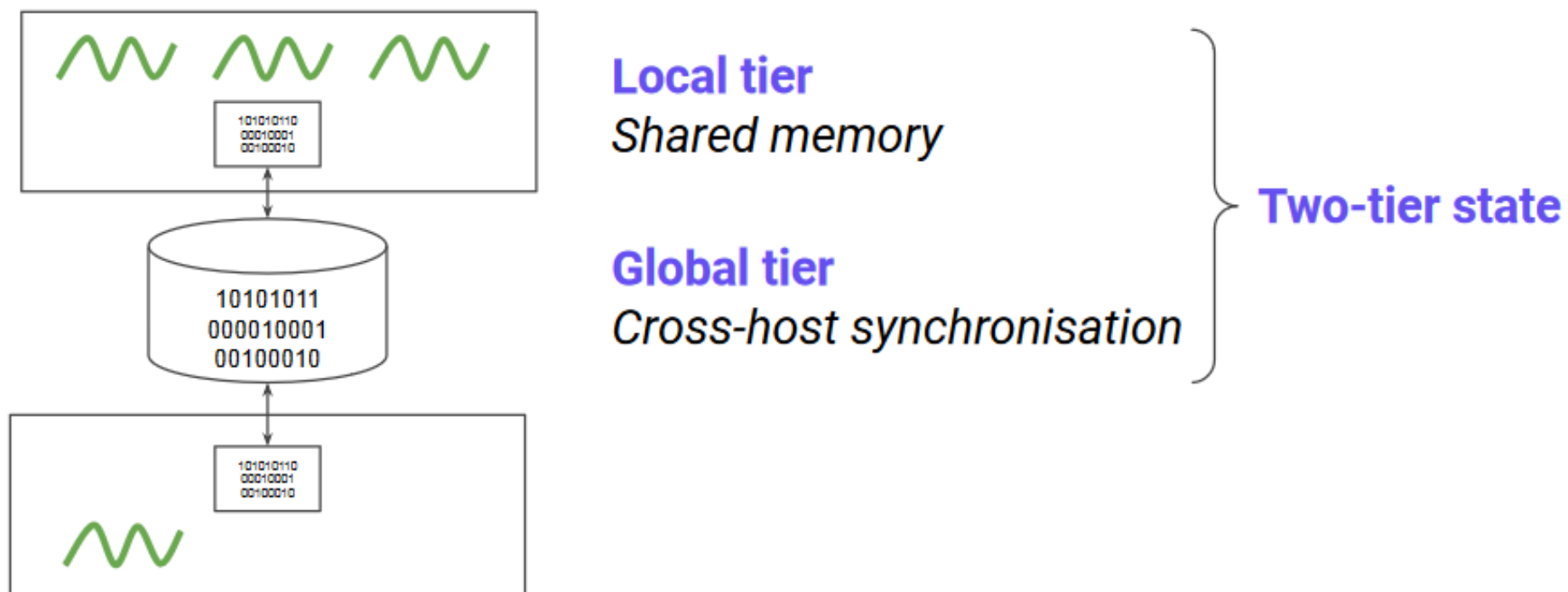
- Lightweight memory safety
- Used by Fastly, Cloudflare, Krustlet



Challenges:

- *Relax isolation to share memory at runtime*
- *Virtualisation between functions and host resources*

Two-Tier State - Distribution *and* Locally-Shared State



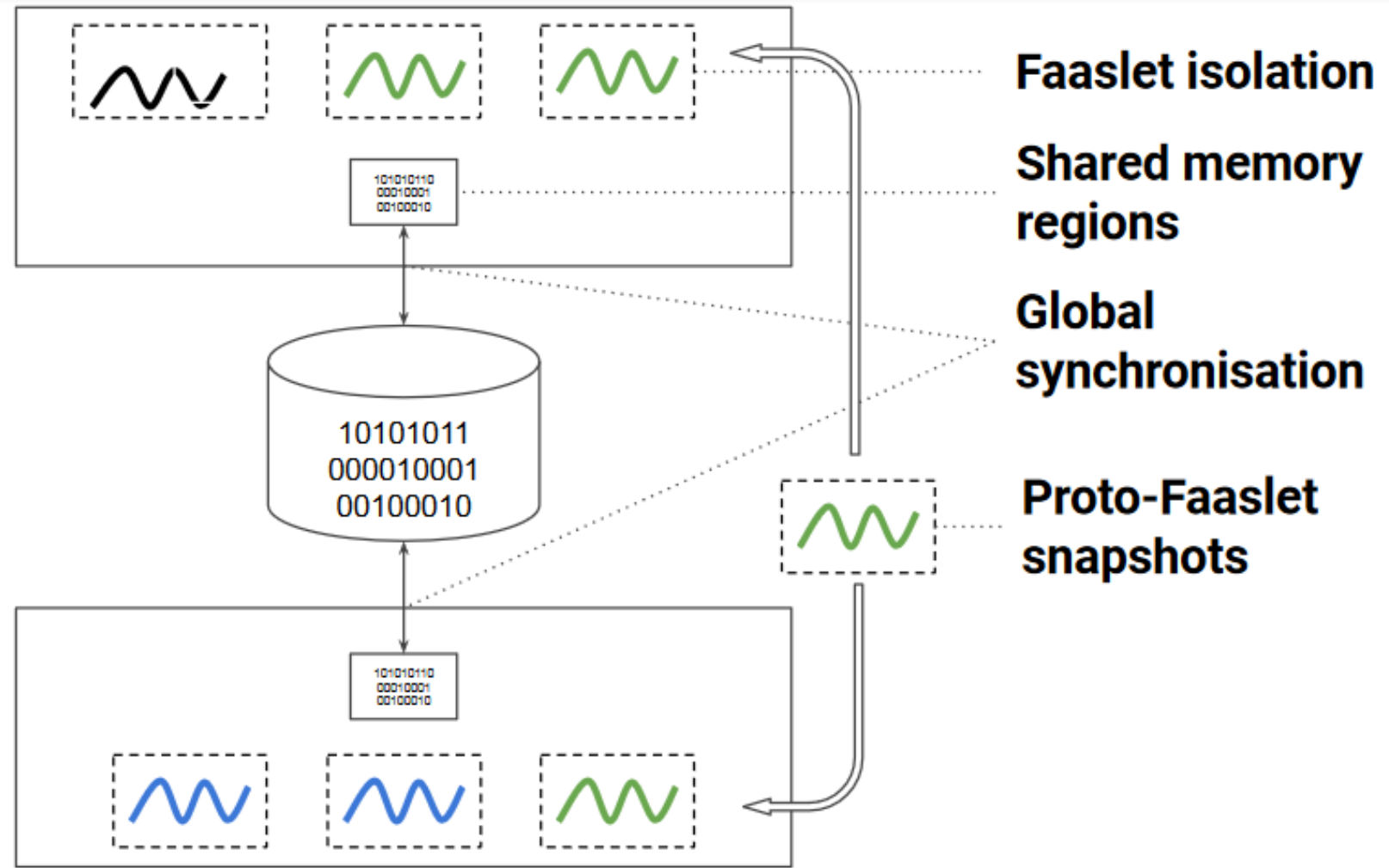
Challenges:

- *Hide complexity from the user*
- *Minimise synchronisation*
- *Schedule to optimise co-location*

Faasm: Lightweight Isolation for Efficient Stateful Serverless Computing



<https://github.com/llds/Faasm>



Problem 1: Isolation overheads

Faaslets - lightweight isolation based on WebAssembly

Host interface - minimal serverless-specific virtualisation

Proto-Faaslets - 500 μ s initialisation, 90kB memory

Problem 2: Inefficient state sharing

Faaslet shared regions - shared memory without breaking isolation

Two-tier state - global synchronisation

Problem 1: Isolation overheads

Faaslets - lightweight isolation based on WebAssembly

Host interface - minimal serverless-specific virtualisation

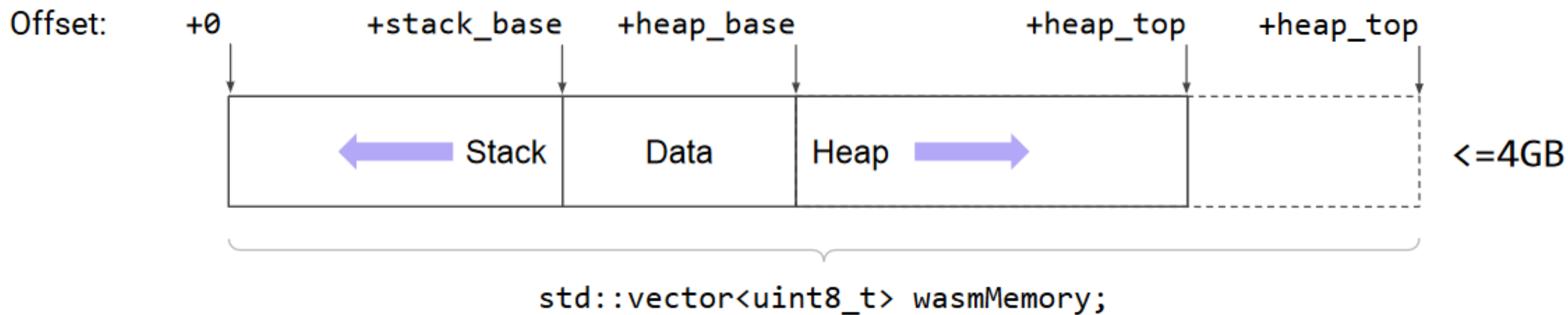
Proto-Faaslets - 500 μ s initialisation, 90kB memory

Problem 2: Inefficient state sharing

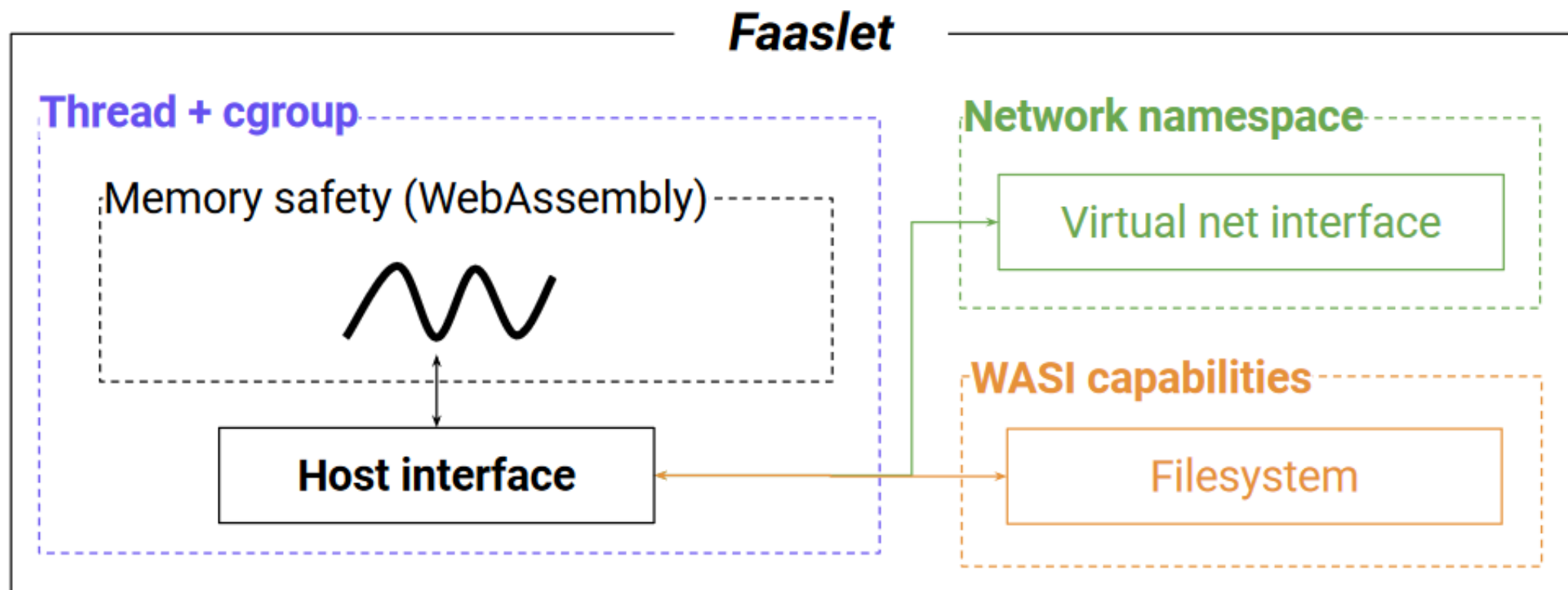
Faaslet shared regions - shared memory without breaking isolation

Two-tier state - global synchronisation

WebAssembly - memory safety with fine-grained control



WebAssembly memory model



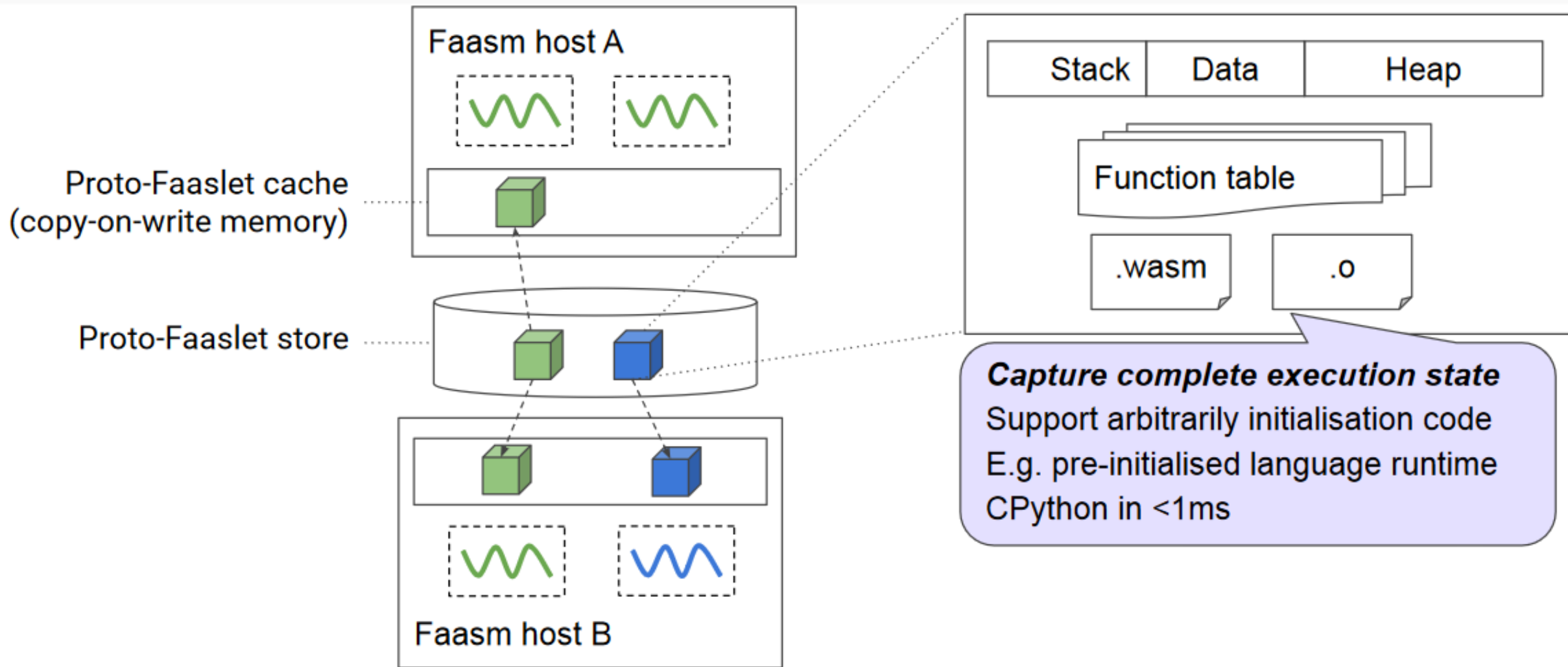
Faaslet multi-tenant isolation

Minimal Virtualisation for Serverless *and* POSIX applications

<i>Category</i>	<i>Sub-category</i>	<i>API</i>
Serverless	Chaining	<code>chain_call()</code> , <code>await_call()</code> , ...
	State	<code>get_state()</code> , <code>set_state()</code> , ...
POSIX	Dynamic Linking	<code>dlopen()</code> , <code>dlsym()</code> , ...
	Memory	<code>mmap()</code> , <code>brk()</code> , ...
	Network	<code>socket()</code> , <code>connect()</code> , <code>bind()</code> , ...
	File I/O	<code>open()</code> , <code>close()</code> , <code>read()</code> , ...

The Faaslet Host Interface

Proto-Faaslets - Host-Independence, μ s Restore, kB Memory Footprint



Proto-Faaslet snapshot and restore

Problem 1: Isolation overheads

Faaslets - lightweight isolation based on WebAssembly

Host interface - minimal serverless-specific virtualisation

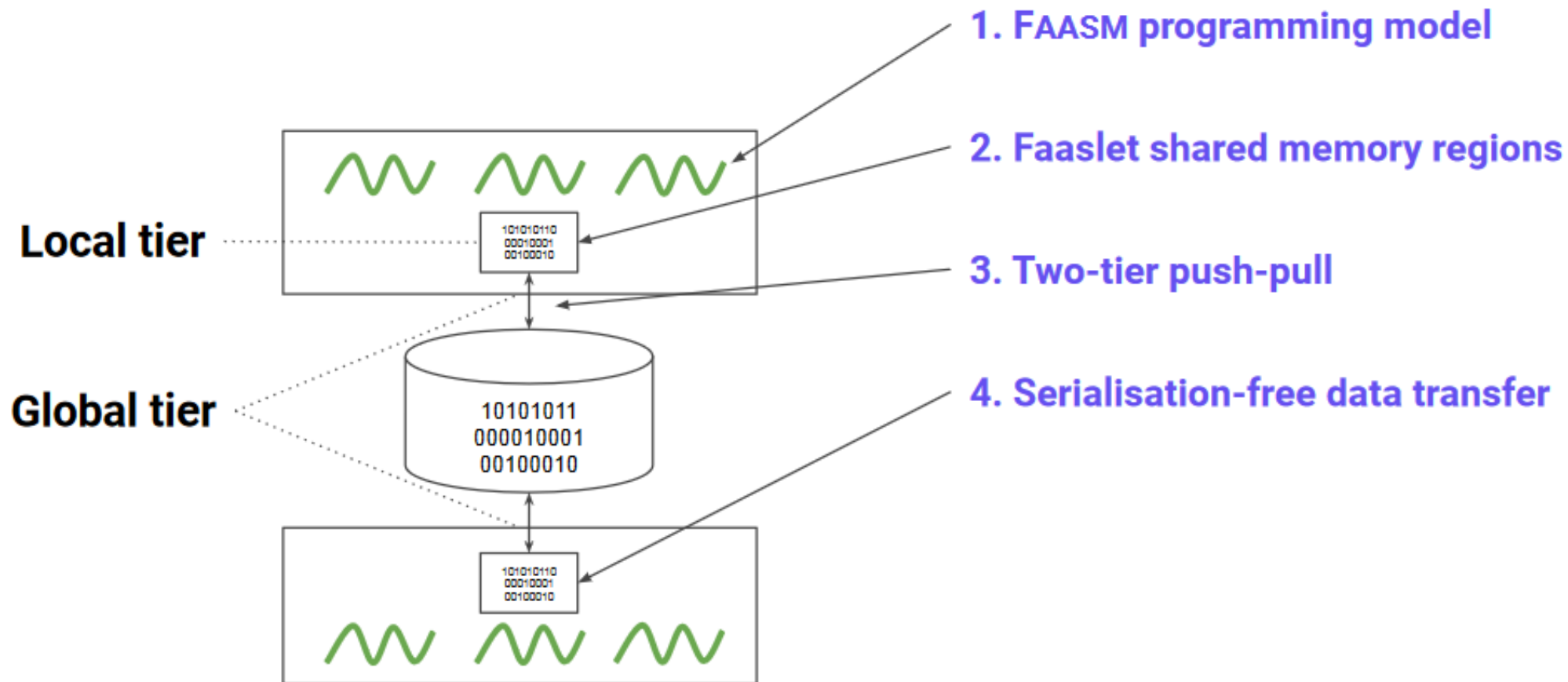
Proto-Faaslets - 500 μ s initialisation, 90kB memory

Problem 2: Inefficient state sharing

Faaslet shared regions - shared memory without breaking isolation

Two-tier state - global synchronisation

Two-Tier State Architecture Top-Down View



FAASM Programming Model - Distributed Machine Learning (SGD)

```
t_a = SparseMatrixReadOnly("training_a")
t_b = MatrixReadOnly("training_b")
weights = VectorAsync("weights")
```

```
@serverless_func
```

```
def weight_update(idx_a , idx_b):
```

```
    for col_idx , col_a in t_a.columns[idx_a:idx_b]:
```

```
        col_b = t_b.columns[col_idx]
```

```
        adj = calc_adjustment(col_a , col_b)
```

```
        for val_idx , val in col_a.non_nulls ():
```

```
            weights[val_idx] += val * adj
```

```
            if iter_count % threshold == 0:
```

```
                weights.push()
```

```
@serverless_func
```

```
def sgd_main(n_workers , n_epochs):
```

```
    for e in n_epochs:
```

```
        args = divide_problem(n_workers)
```

```
        c = chain(weight_update, n_workers, args)
```

```
    await_all(c)
```

High-level Object-Oriented abstractions

Read-only matrices

Asynchronous vector

Flexible consistency

Standard Programming constructs

Transparent optimisations

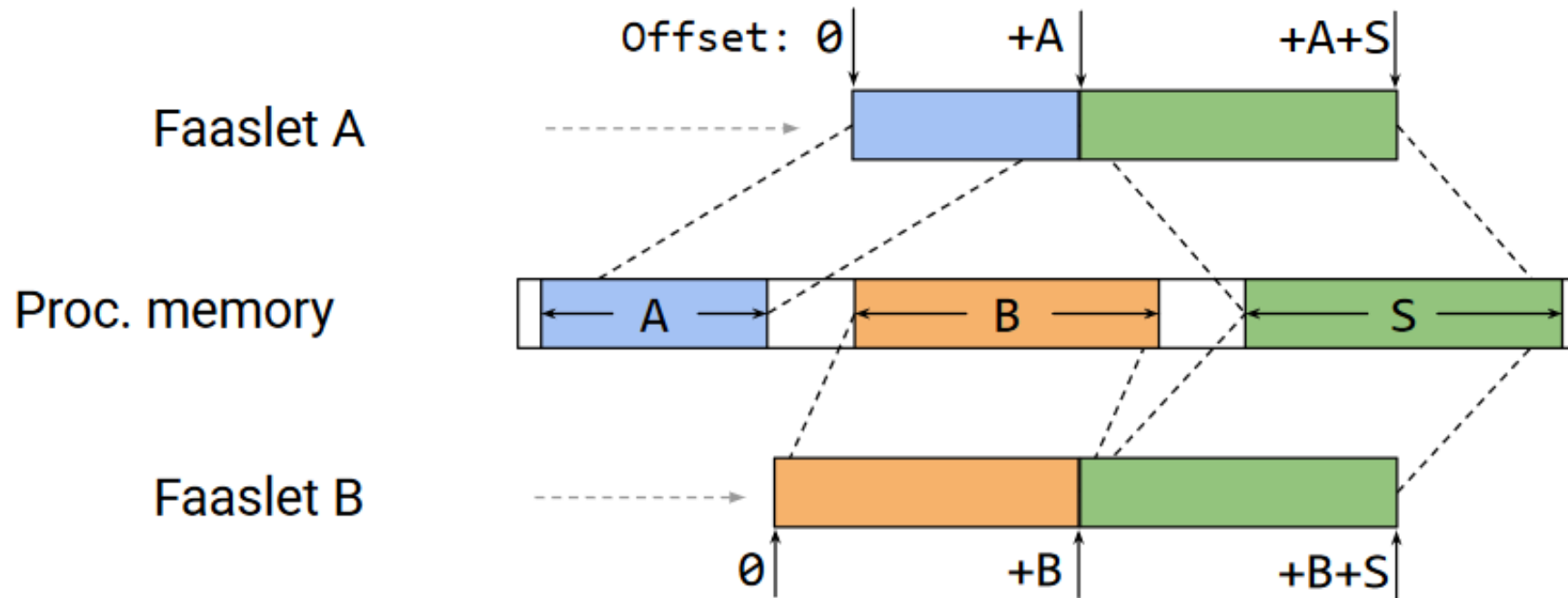
Direct access to shared memory

Intuitive mark-up

Function annotation

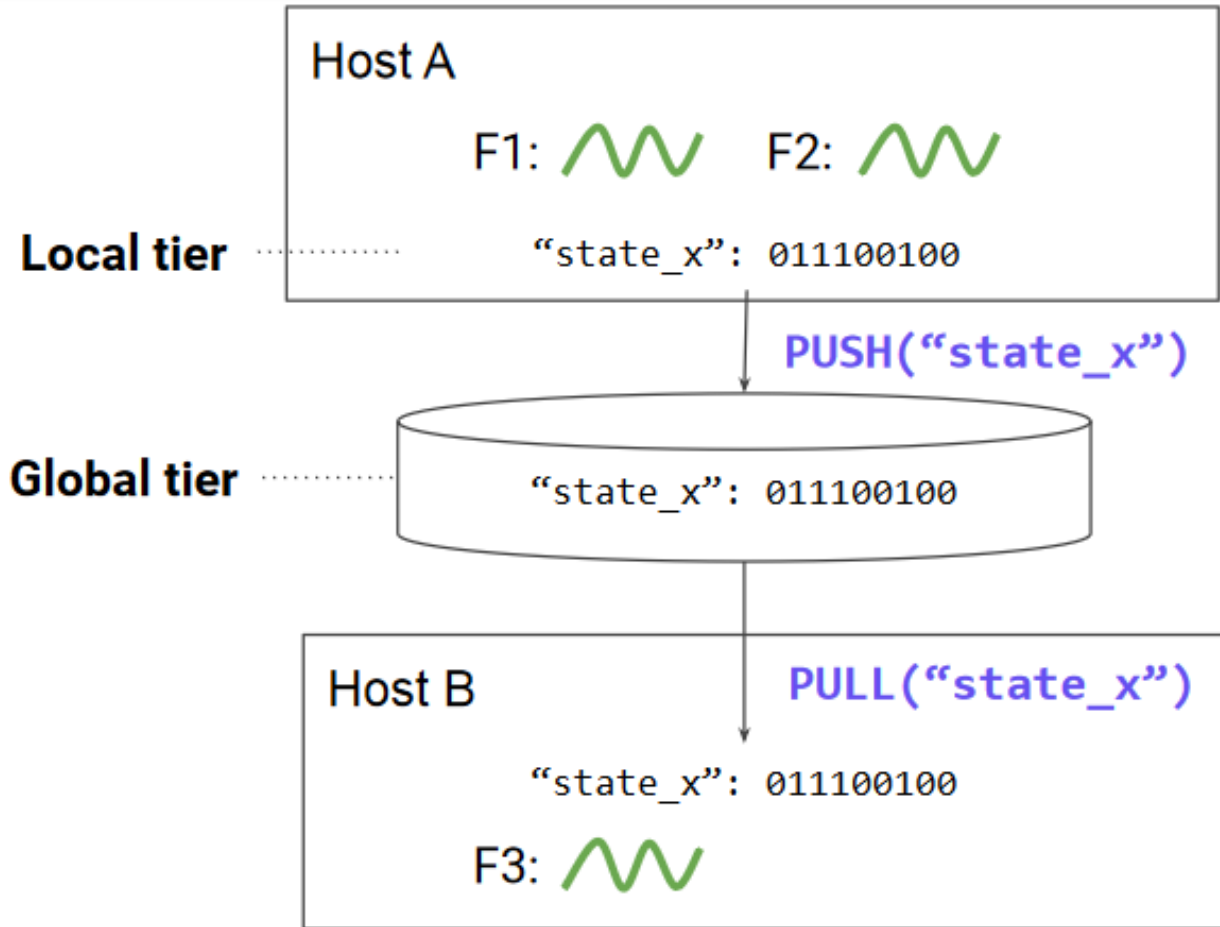
Fork-join parallelism

Shared Memory Without Breaking Safety Guarantees



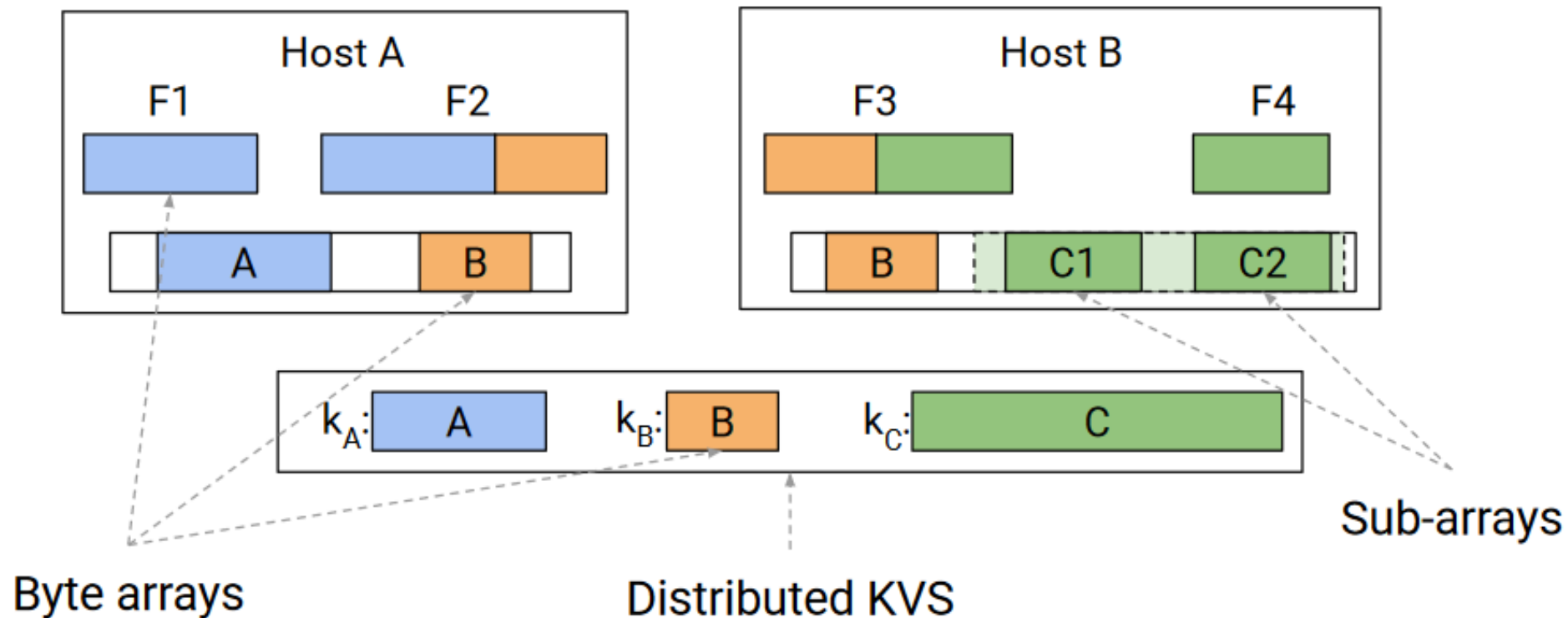
Faaslet Shared Memory Regions

Push-pull - Global Synchronisation with Variable Consistency



Two-Tier Push-Pull

Serialisation-Free Transfer of Arbitrarily Complex Data Structures



Faasm's serialisation-free state

Questions:

1. How do Faaslets compare to containers?
2. Can FAASM improve efficiency and performance of ML training?
3. Can FAASM improve throughput of ML inference?
4. Does Faaslet isolation affect performance of dynamic languages?

Comparison:

- Knative running identical code
- Code compiled natively for Knative
- Code compiled to WebAssembly for FAASM



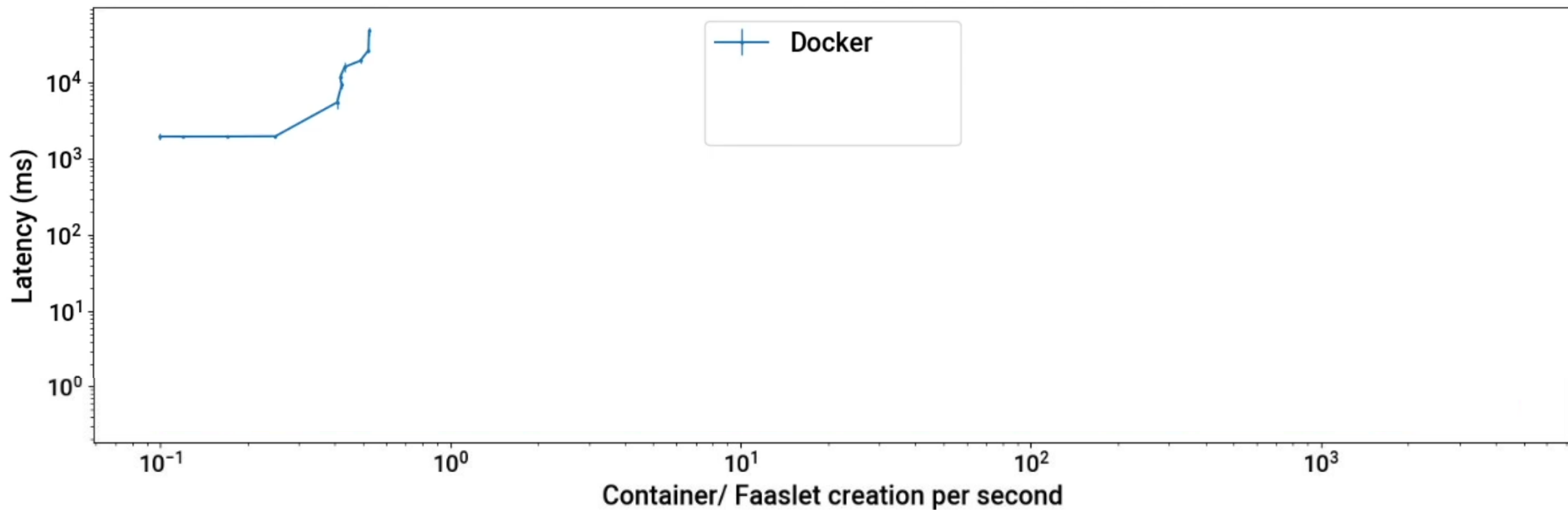
How do Faaslet Overheads Compare to Containers?

Lower overheads mean lower latency and lower costs

	Docker (alpine)	Faaslets	Proto-Faaslets	vs. Docker
Initialisation	2.8s	5.2ms	0.5ms	5.6K x
CPU cycles	251M	1.4K	650	385K x
Memory Footprint	1.3MB	200KB	90KB	15 x
Density	~8K	~70K	>100K	12 x

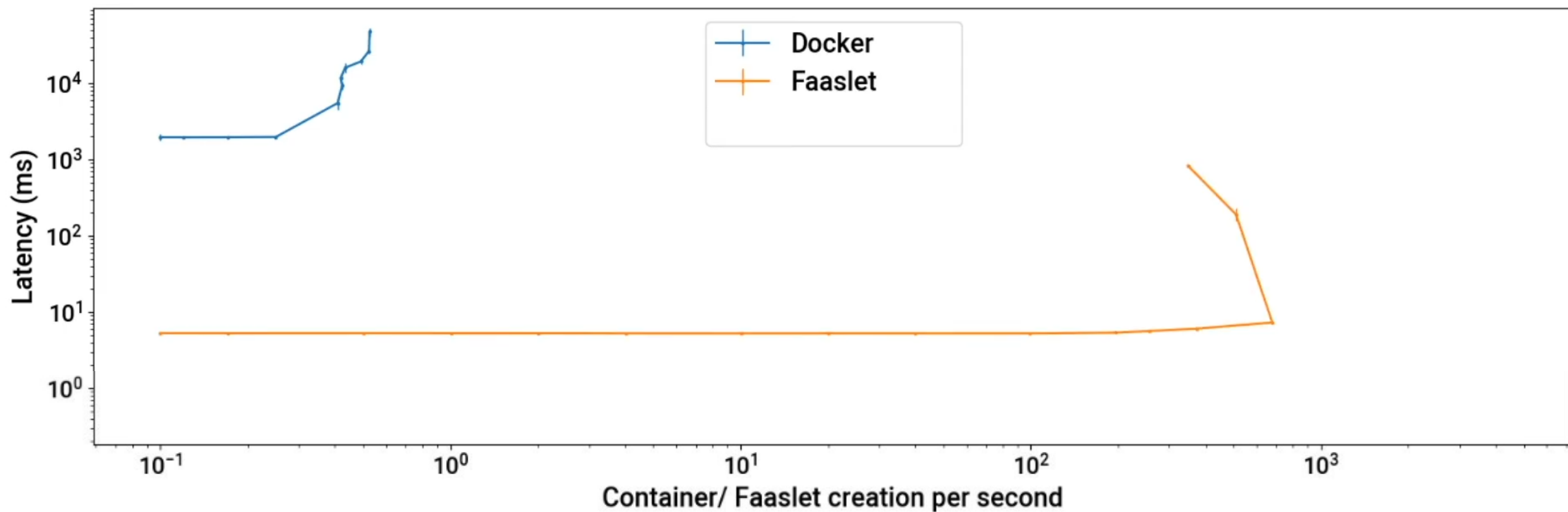
How do Faaslets "Churn" Compared to Containers?

Higher churn means higher utilisation of shared infrastructure



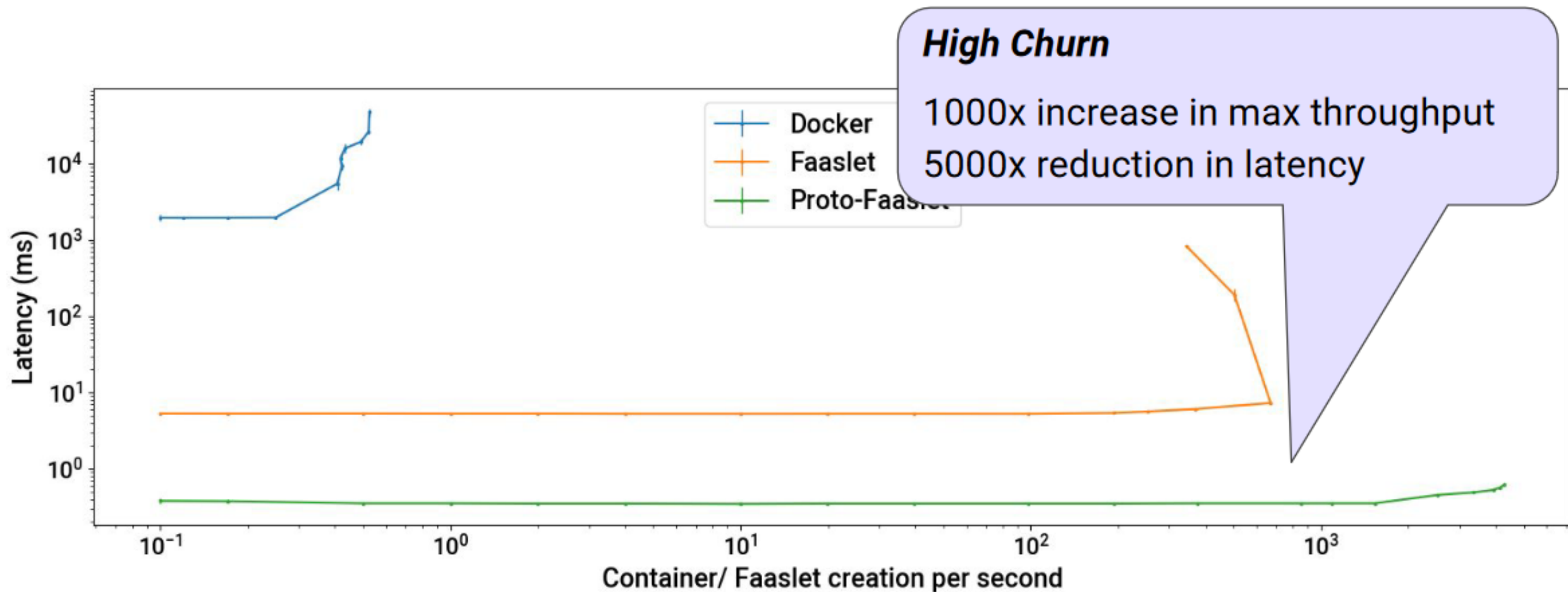
How do Faaslets "Churn" Compared to Containers?

Higher churn means higher utilisation of shared infrastructure

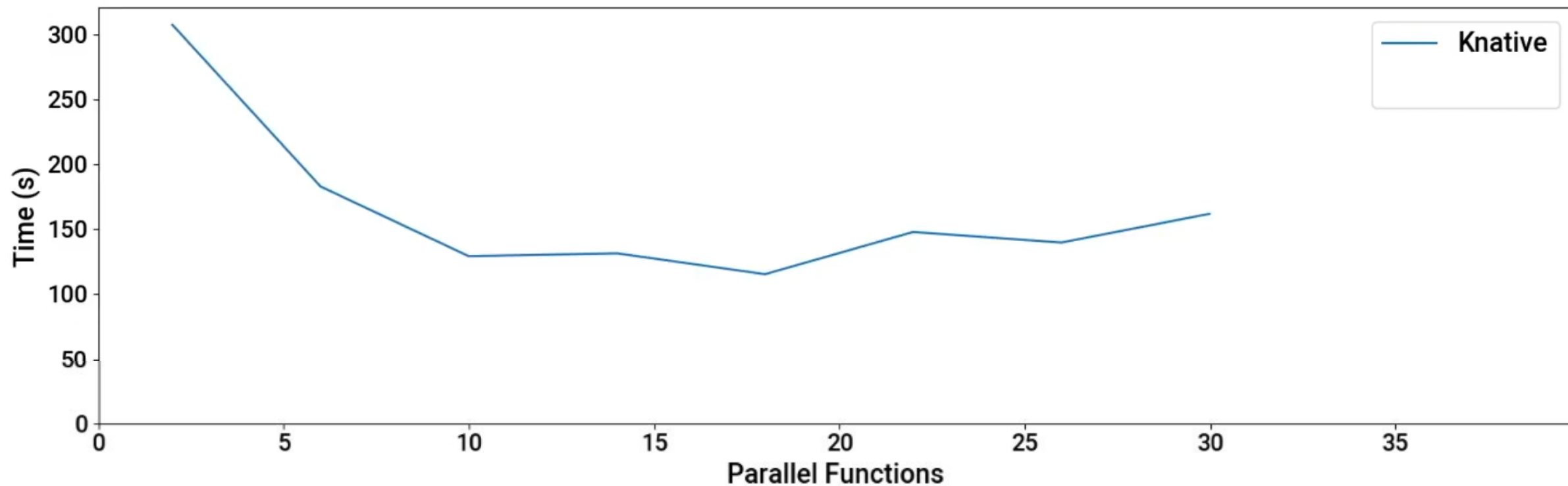


How do Faaslets "Churn" Compared to Containers?

Higher churn means higher utilisation of shared infrastructure

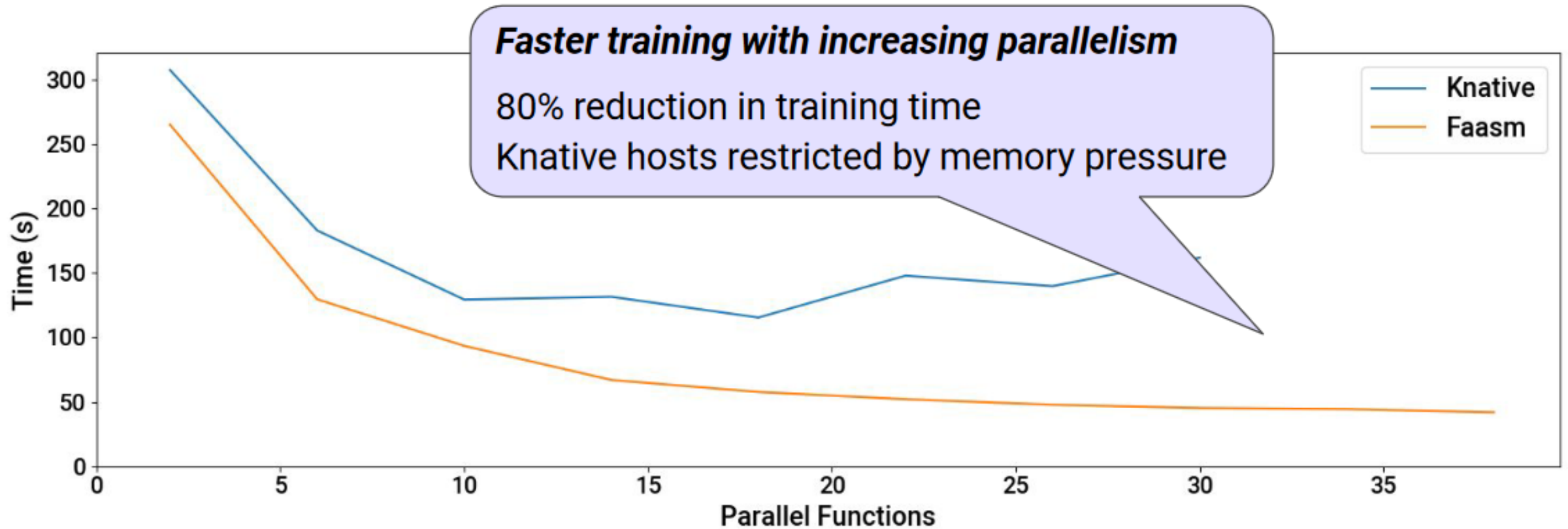


Parallel processing on co-located data reduces training time

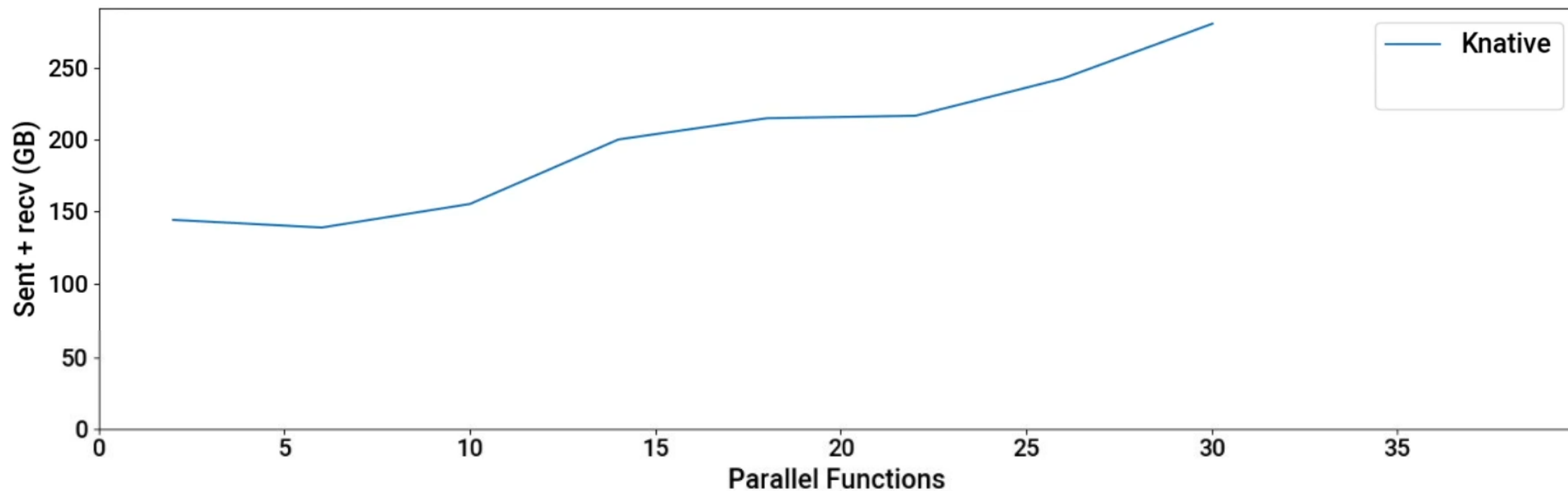


Can Faasm Improve Efficiency and Performance of Parallel ML Training?

Parallel processing on co-located data reduces training time



Reduced data shipping reduces costs



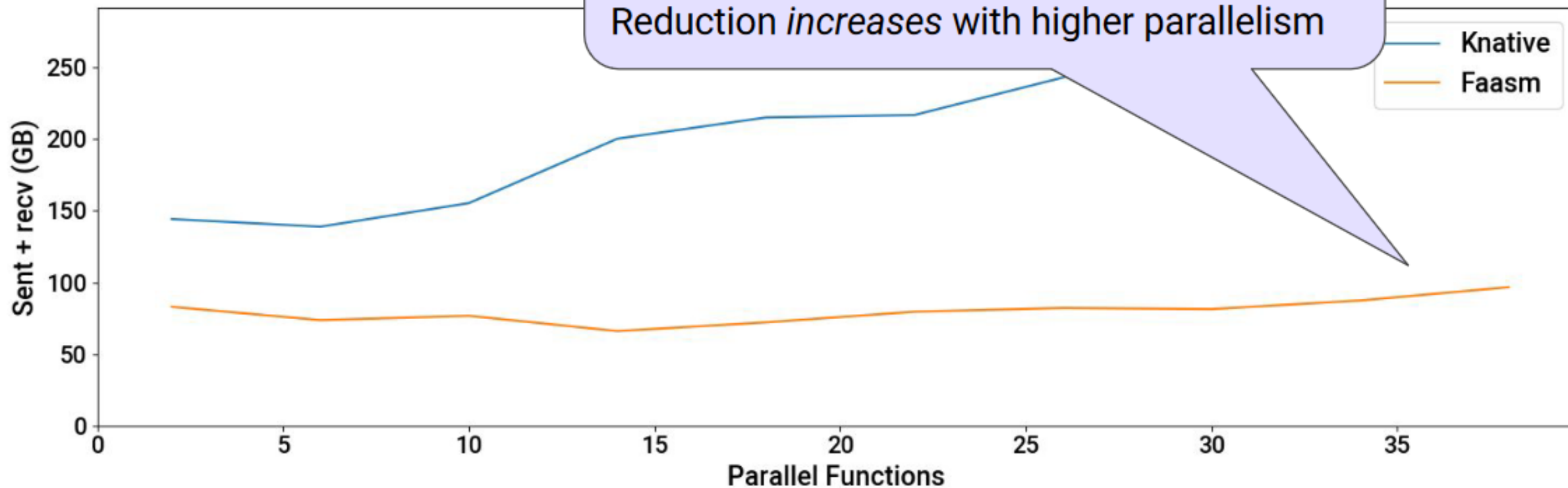
Can Faasm Improve Efficiency and Performance of Parallel ML Training?

Reduced data shipping reduces costs

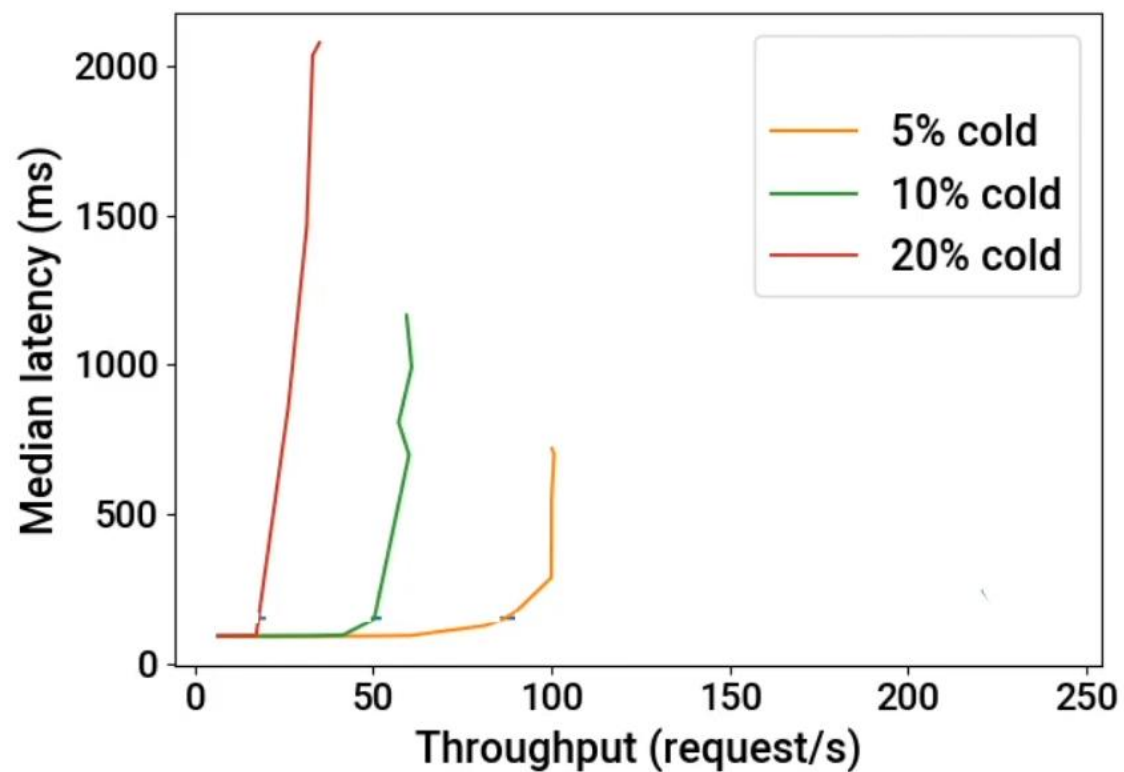
Reduced network transfers

60% reduction in network transfers

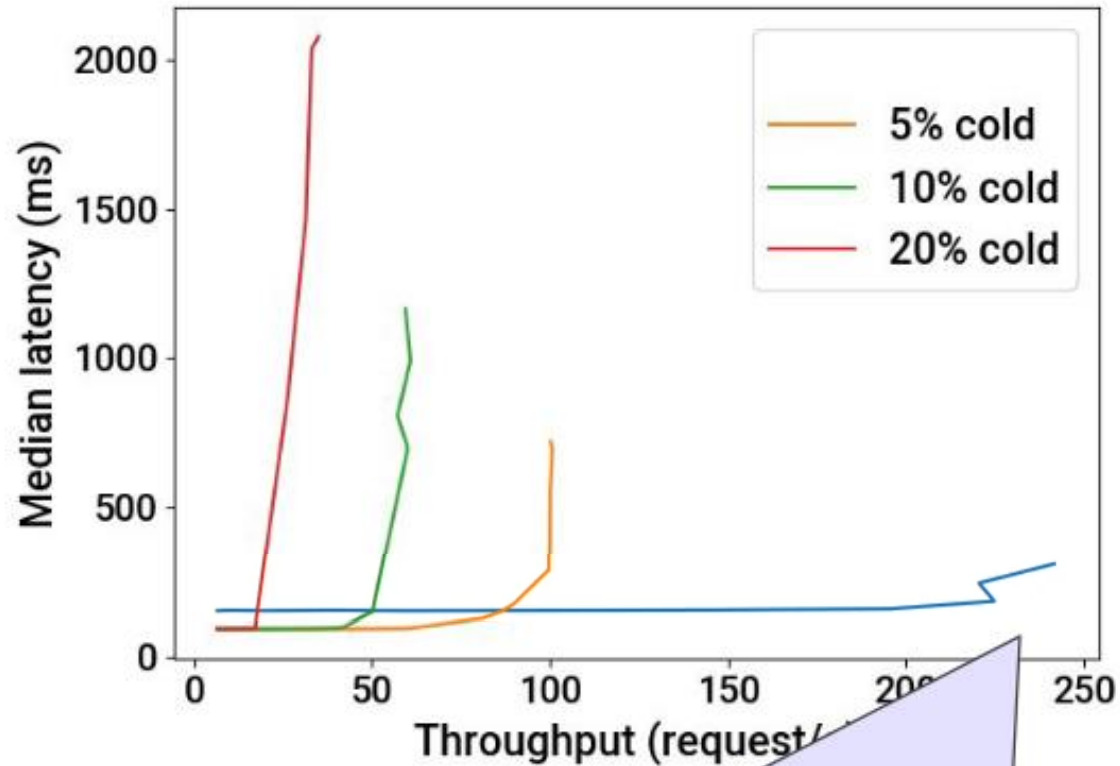
Reduction *increases* with higher parallelism



Proto-Faaslets increase max throughput and reduce latency



Proto-Faaslets increase max throughput and reduce latency



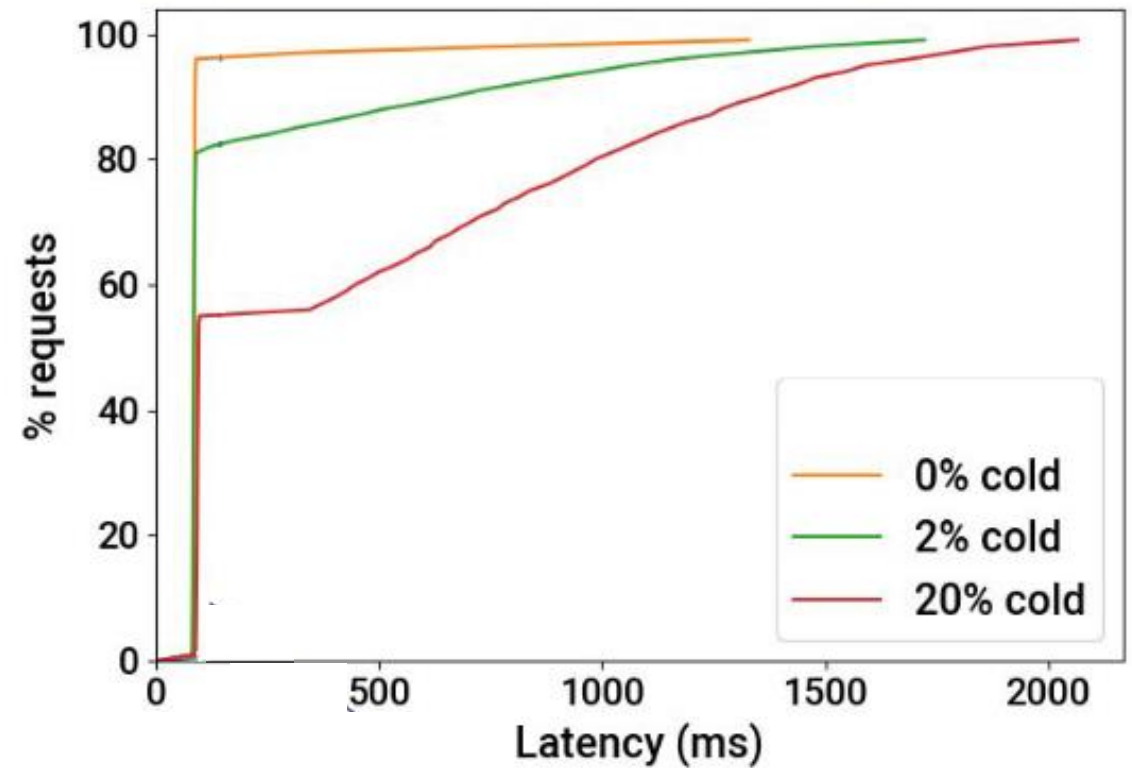
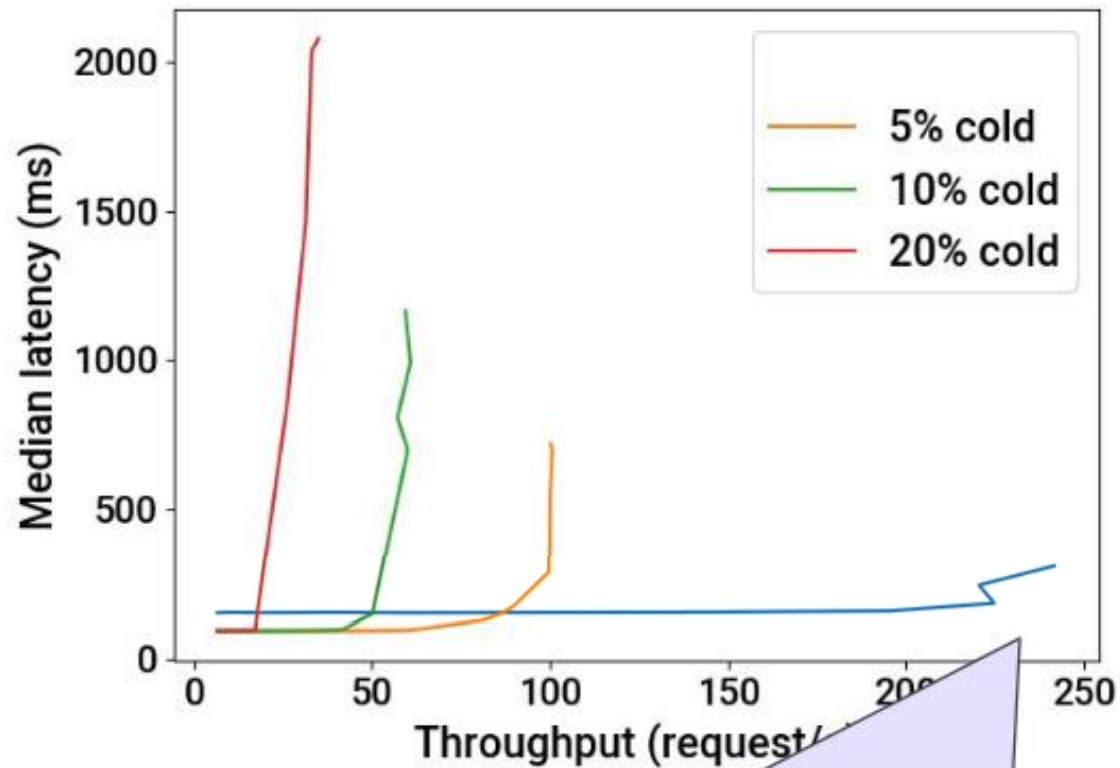
Increased Throughput

Negligible cold starts with Proto-Faaslets

120% increase in max throughput with 5% cold starts

Can Faasm Improve Throughput and Reduce Latency Serving ML Inference?

Proto-Faaslets increase max throughput and reduce latency

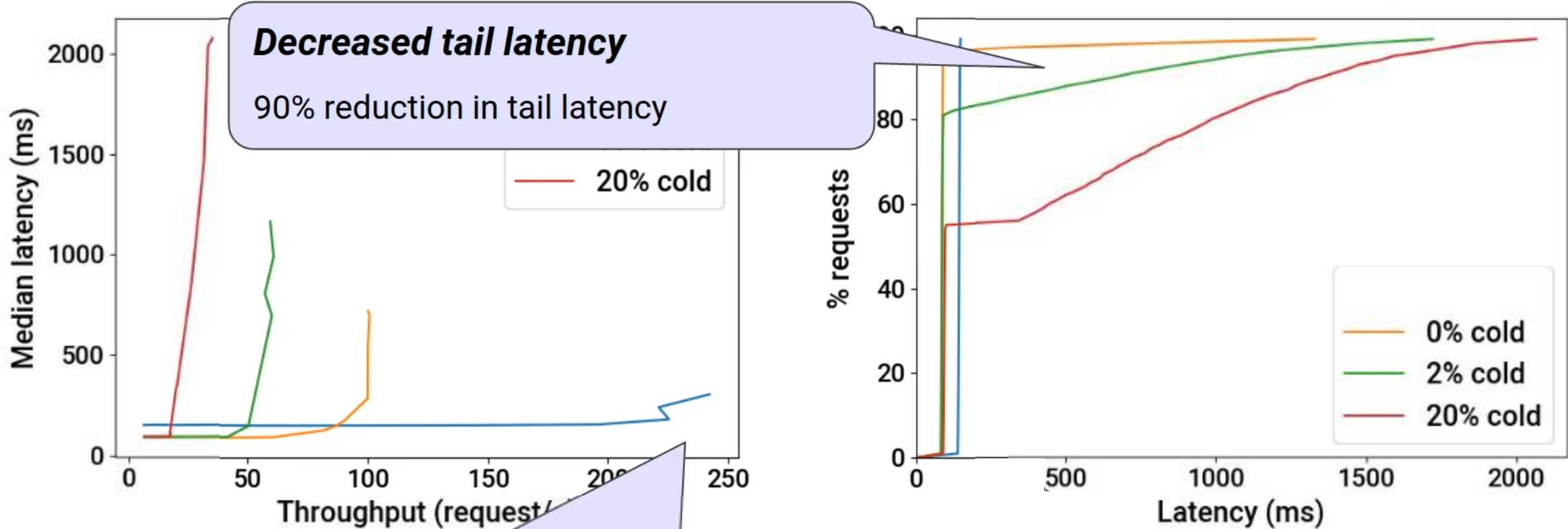


Increased Throughput

Negligible cold starts with Proto-Faaslets
120% increase in max throughput with 5% cold starts

Can Faasm Improve Throughput and Reduce Latency Serving ML Inference?

Proto-Faaslets increase max throughput and reduce latency

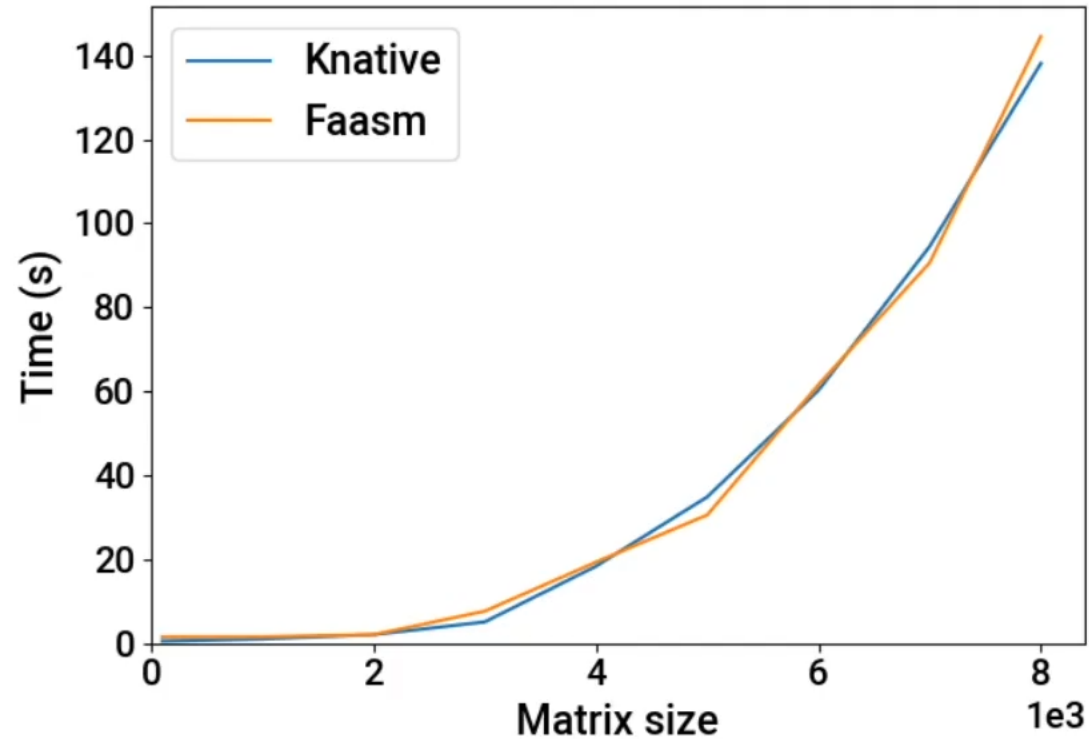


Increased Throughput

Negligible cold starts with Proto-Faaslets
120% increase in max throughput with 5% cold starts

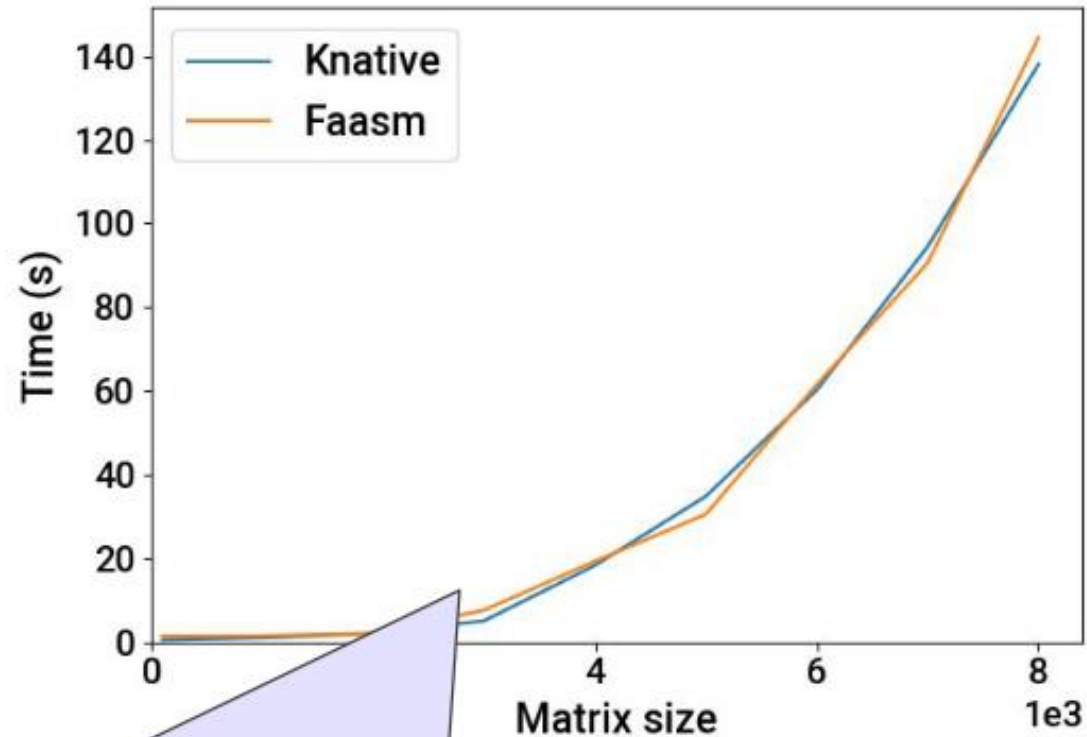
Does Faaslet Isolation Affect Performance of Dynamic Languages?

Faaslet isolation has negligible impact on a distributed Python application



Does Faaslet Isolation Affect Performance of Dynamic Languages?

Faaslet isolation has negligible impact on a distributed Python application

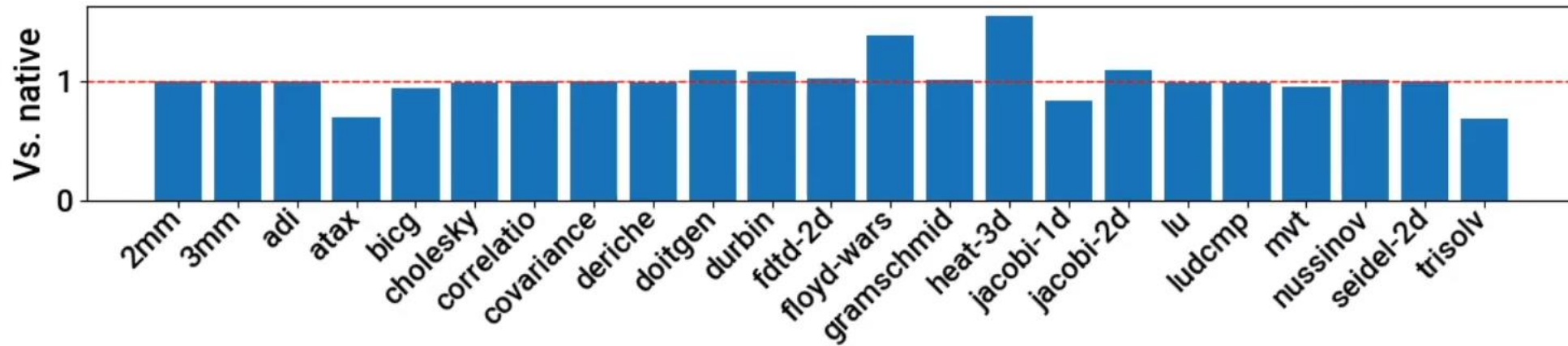


Comparable performance

Faaslet isolation shows no significant overhead
Effect persists with increasing matrix size

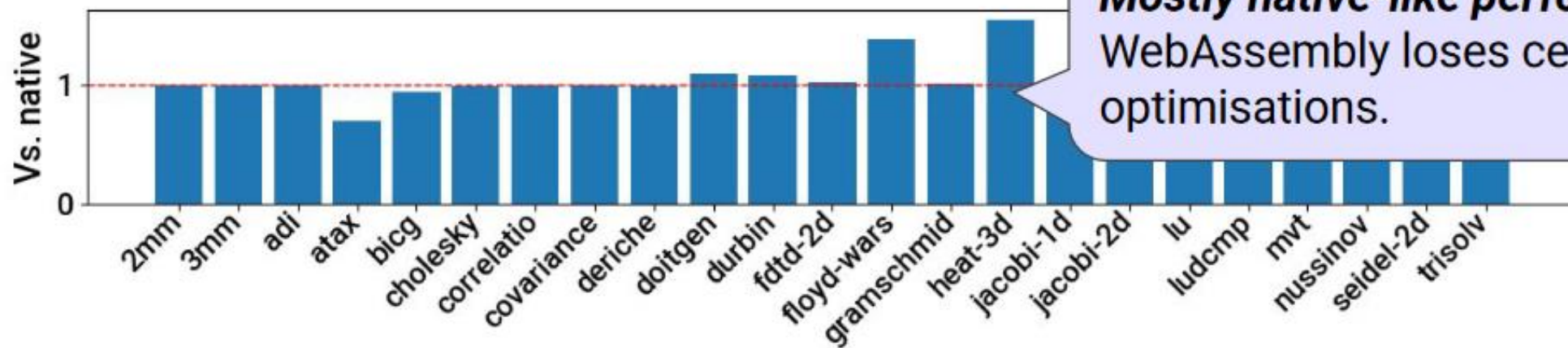
Does Faaslet Isolation Affect Performance of Dynamic Languages?

Performance overheads increase as applications become more complex



Does Faaslet Isolation Affect Performance of Dynamic Languages?

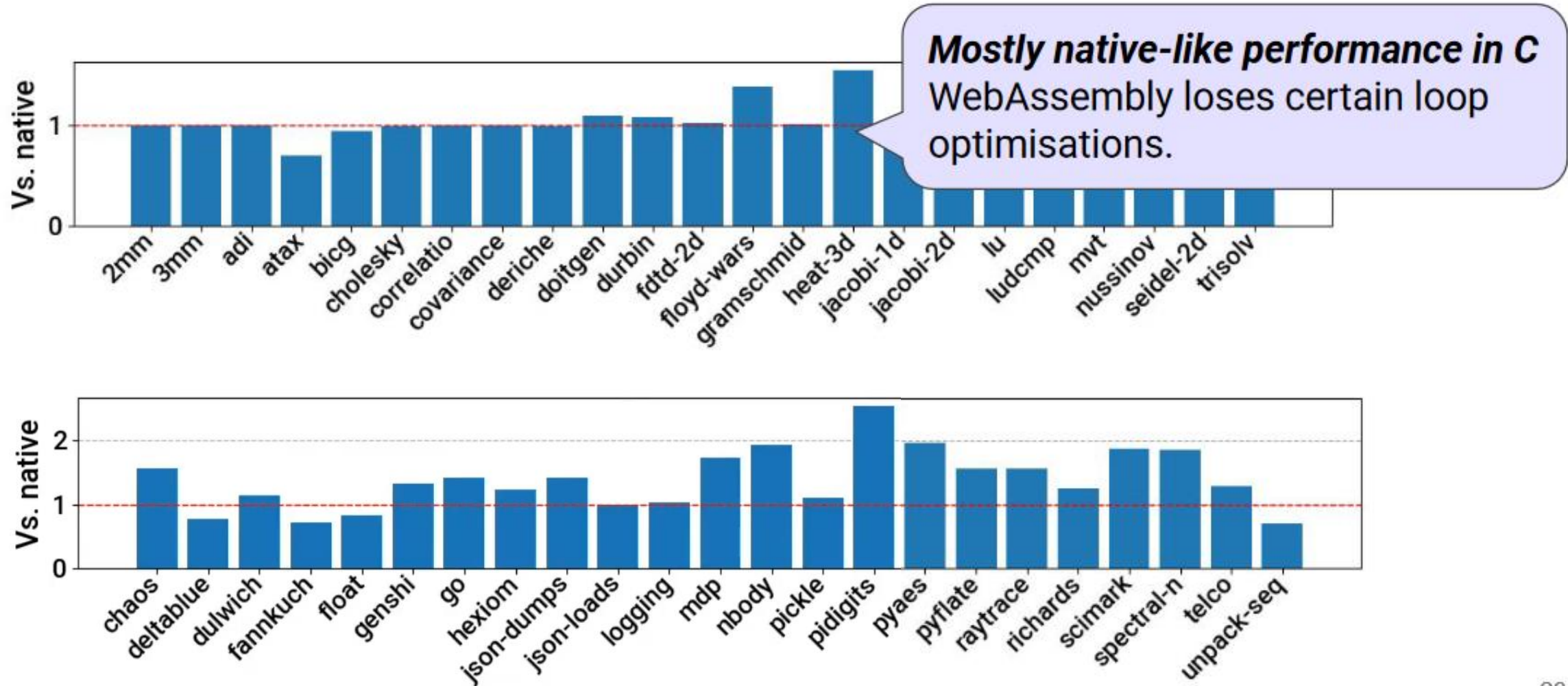
Performance overheads increase as applications become more complex



Mostly native-like performance in C
WebAssembly loses certain loop optimisations.

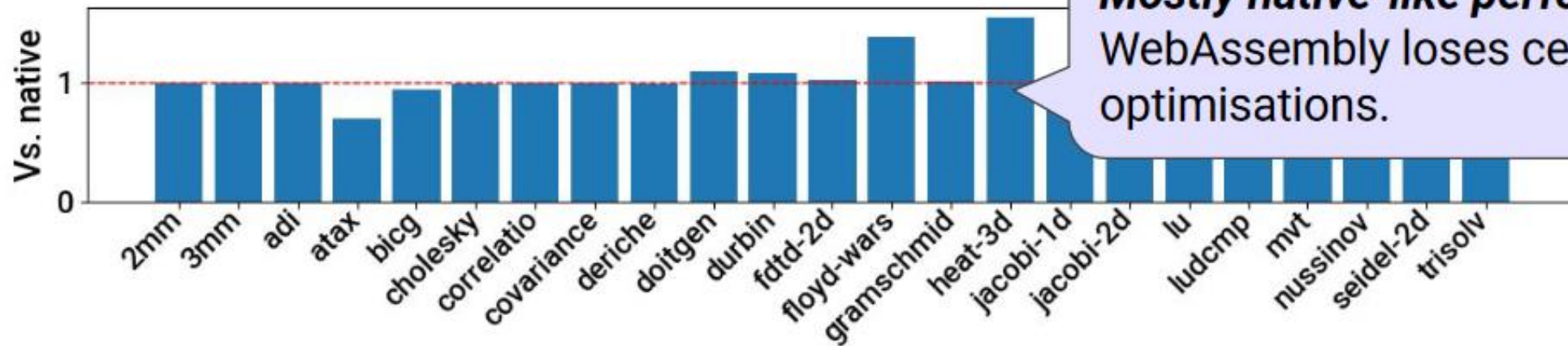
Does Faaslet Isolation Affect Performance of Dynamic Languages?

Performance overheads increase as applications become more complex

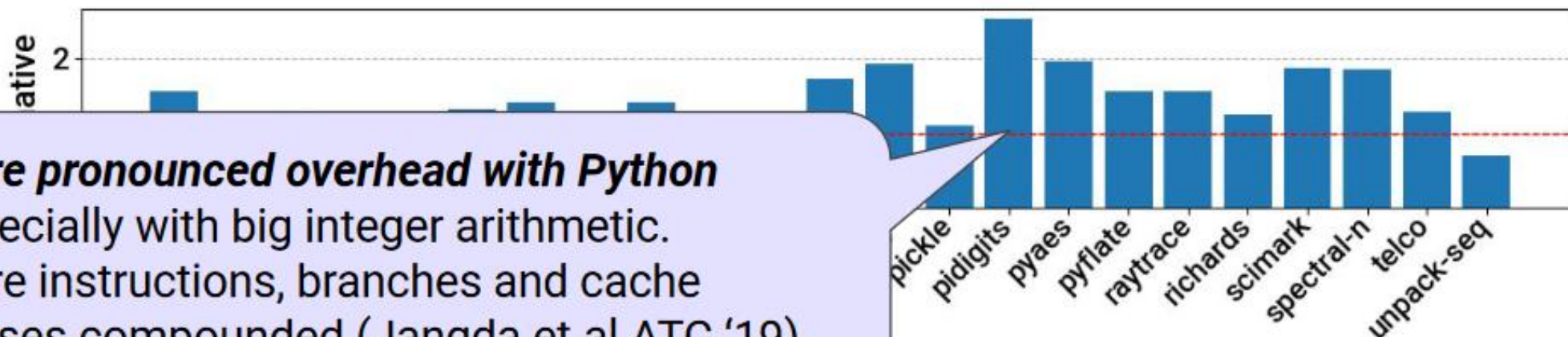


Does Faaslet Isolation Affect Performance of Dynamic Languages?

Performance overheads increase as applications become more complex



Mostly native-like performance in C
WebAssembly loses certain loop optimisations.



More pronounced overhead with Python
Especially with big integer arithmetic.
More instructions, branches and cache misses compounded (Jangda et.al ATC '19).

FAASM makes serverless faster and cheaper:

- Current systems exhibit isolation overhead and inefficient state sharing
- FAASM reduces overheads with Faaslets and Proto-Faaslets
- FAASM supports efficient locally shared and globally synchronised state
- Future work: serverless HPC, trusted hardware, unikernel-based runtime



<https://github.com/llds/Faasm>

Thank you