

# Fast, Distributed Computations in the Cloud

Omid Mashayekhi

Advisor: Philip Levis

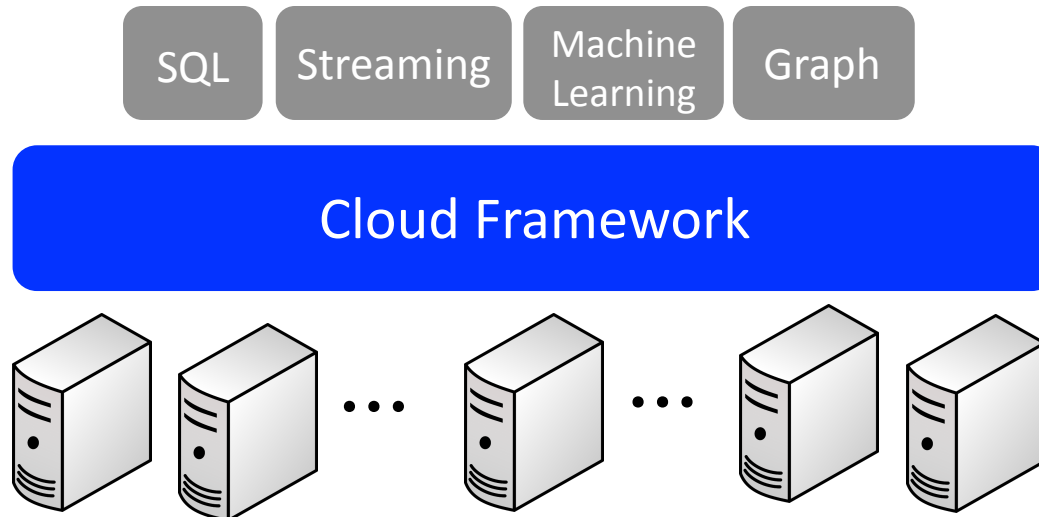


April 7, 2017





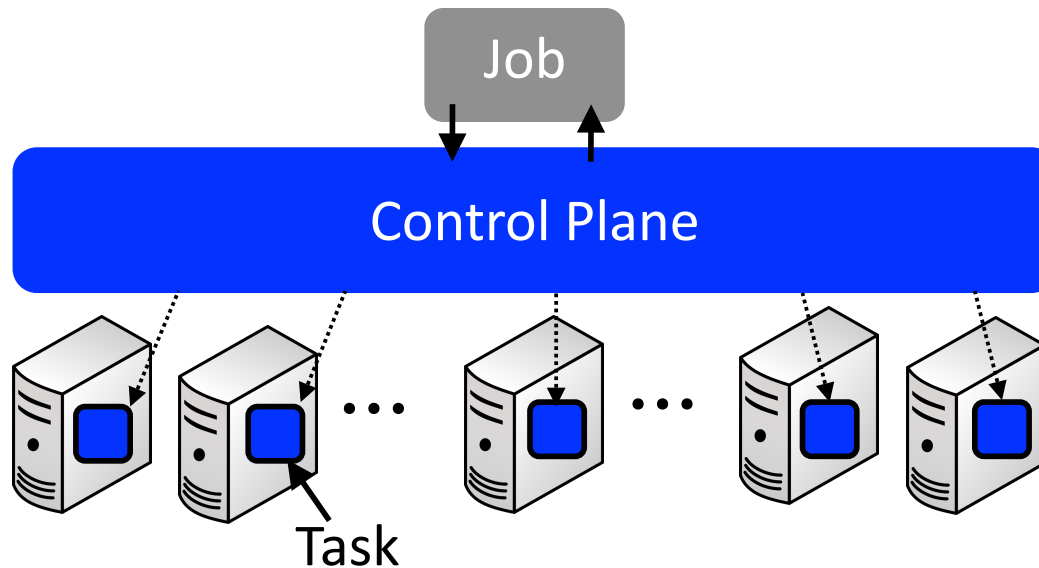
# Cloud Frameworks



Cloud frameworks abstract away the complexities of the cloud infrastructure from the application developers:

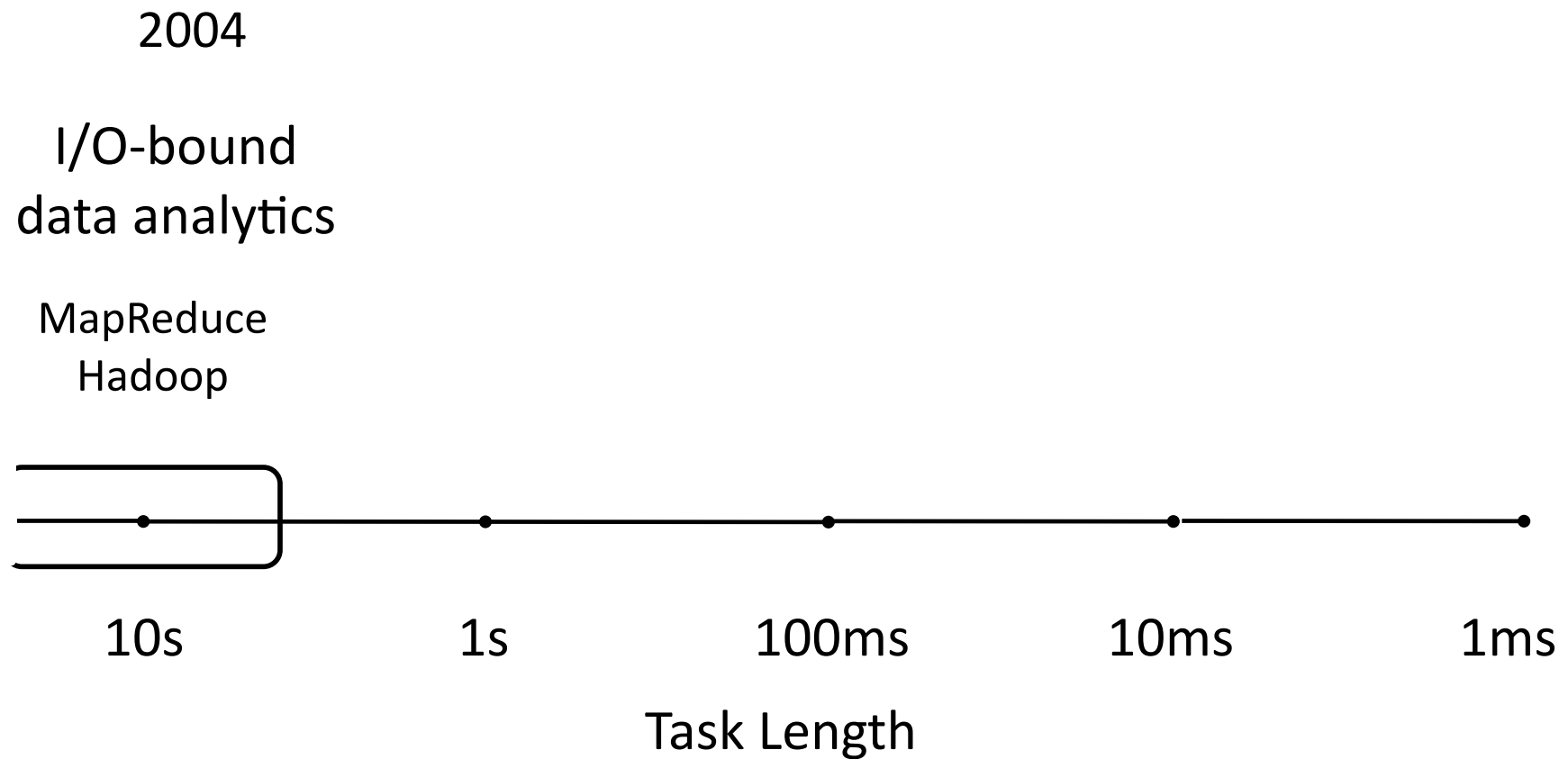
1. Automatic distribution
2. Elastic scalability
3. Multitenant applications
4. Load balancing
5. Fault tolerance

# Cloud Frameworks

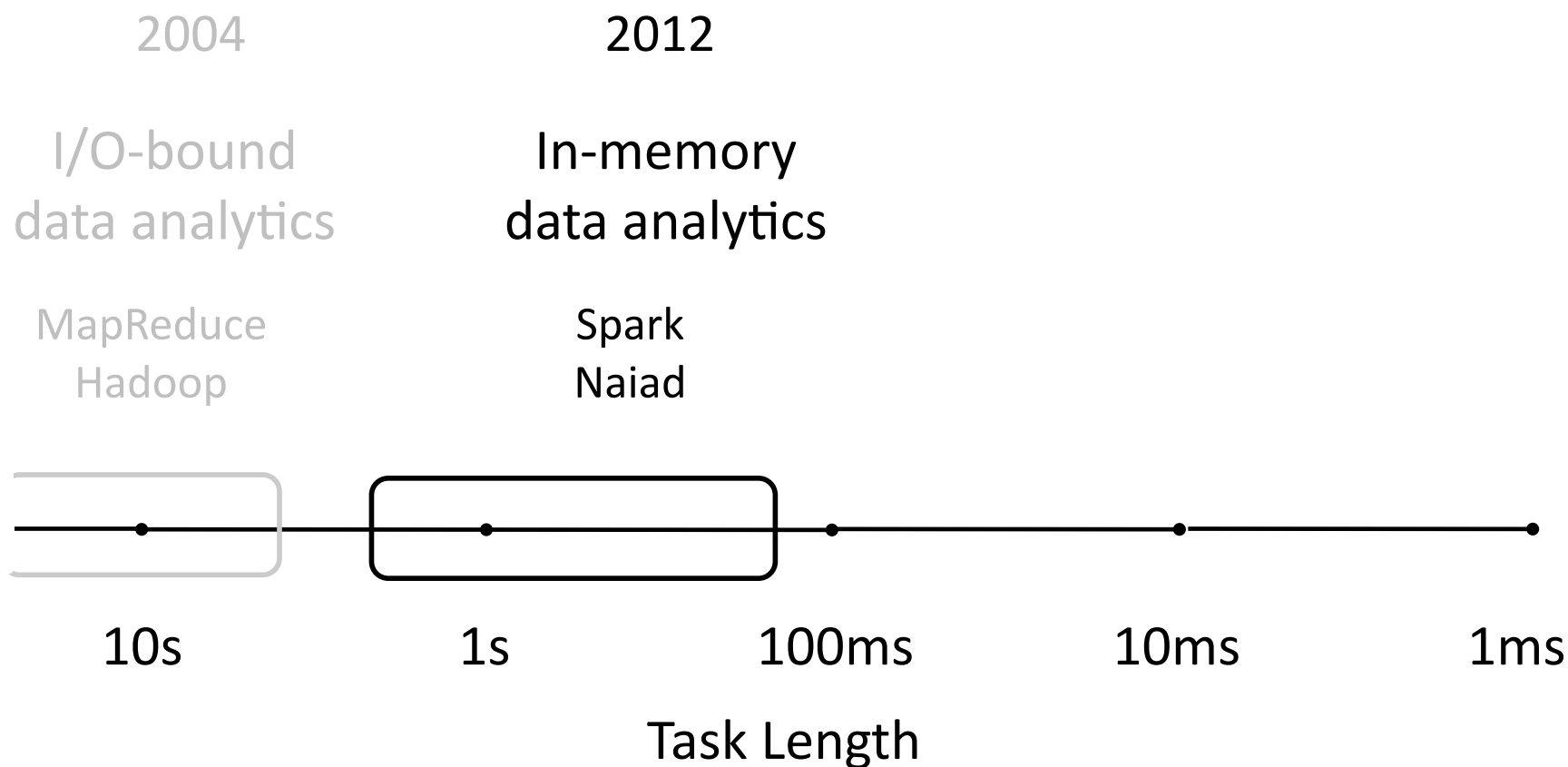


- **Job** is an instance of the application running in the framework.
- **Task** is the unit of computation.
- **Control plane** makes the magic happen:
  - Partitioning job in to tasks
  - Scheduling tasks
  - Load balancing
  - Fault recovery

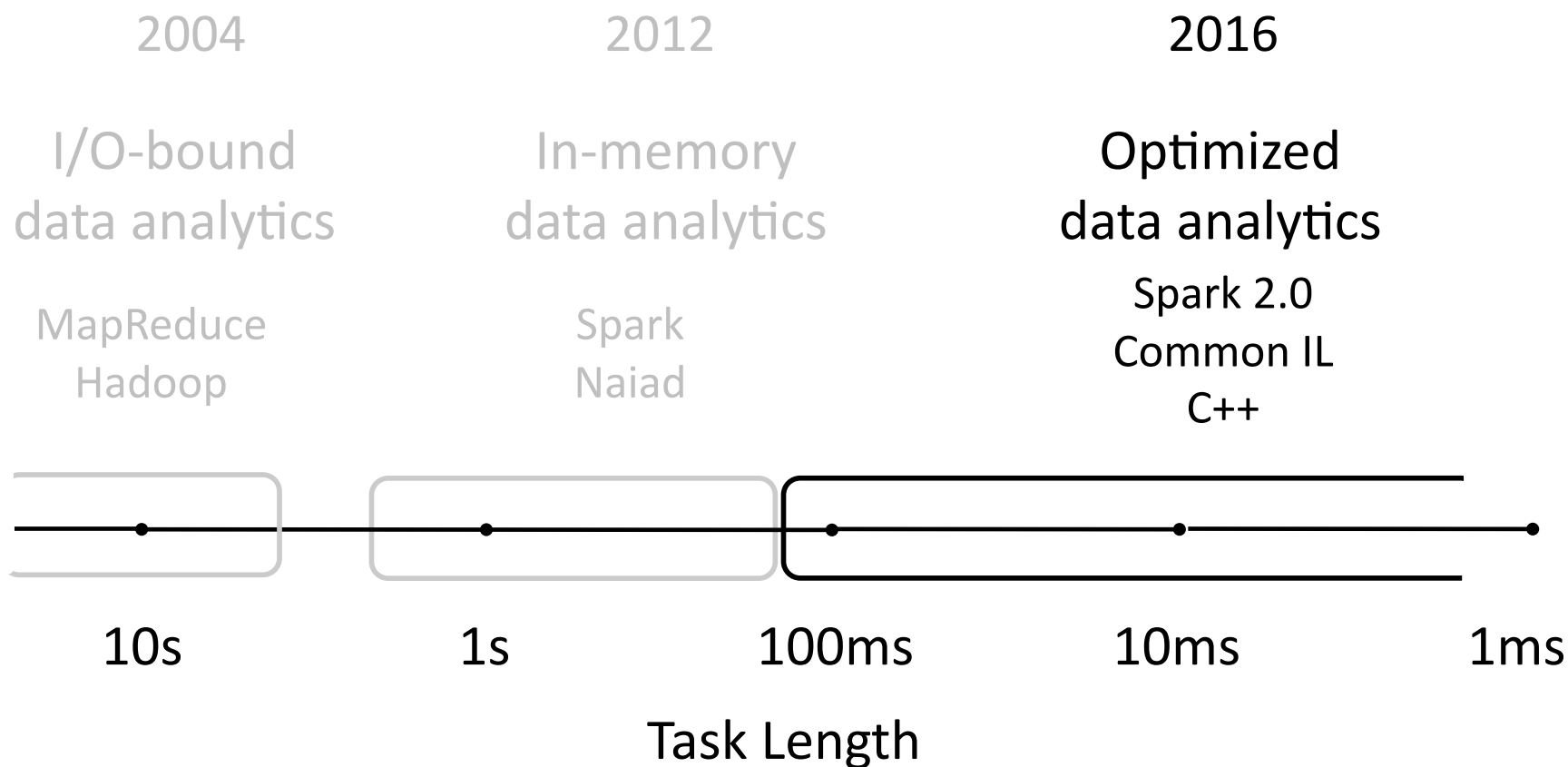
# Evolution of Cloud Frameworks

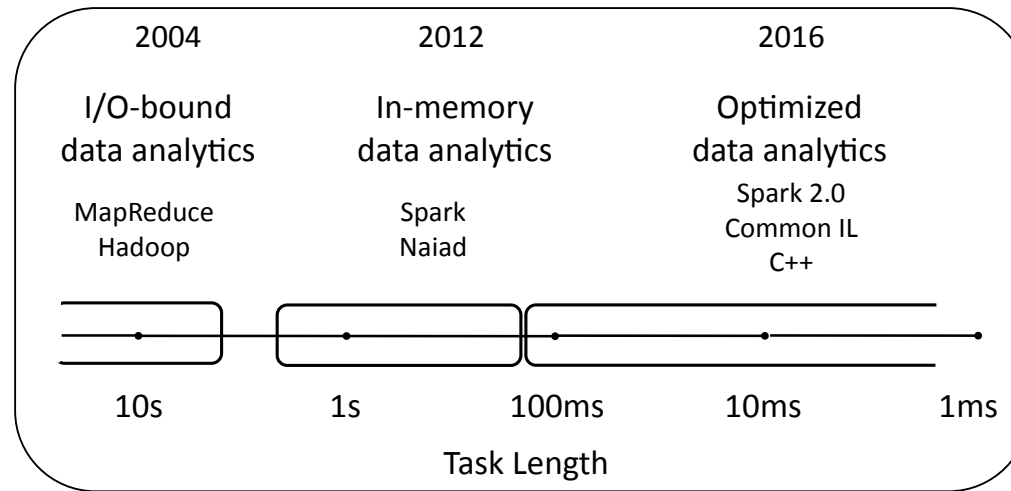


# Evolution of Cloud Frameworks

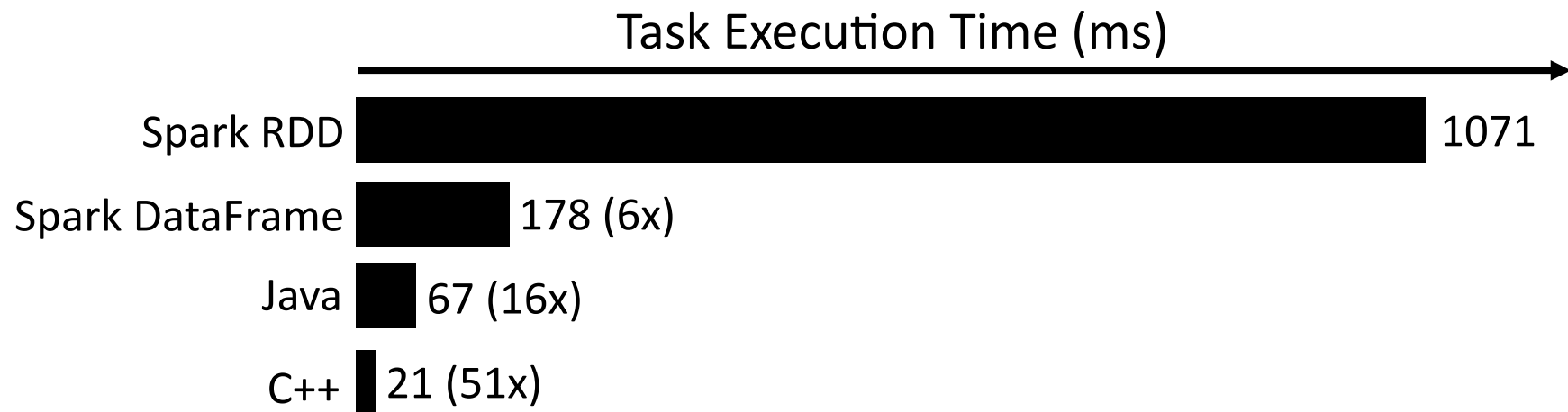
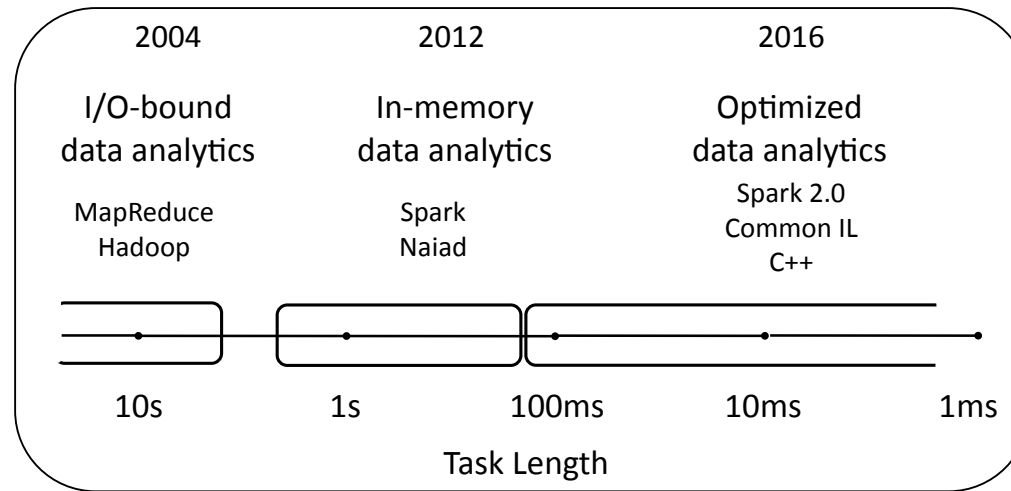


# Evolution of Cloud Frameworks







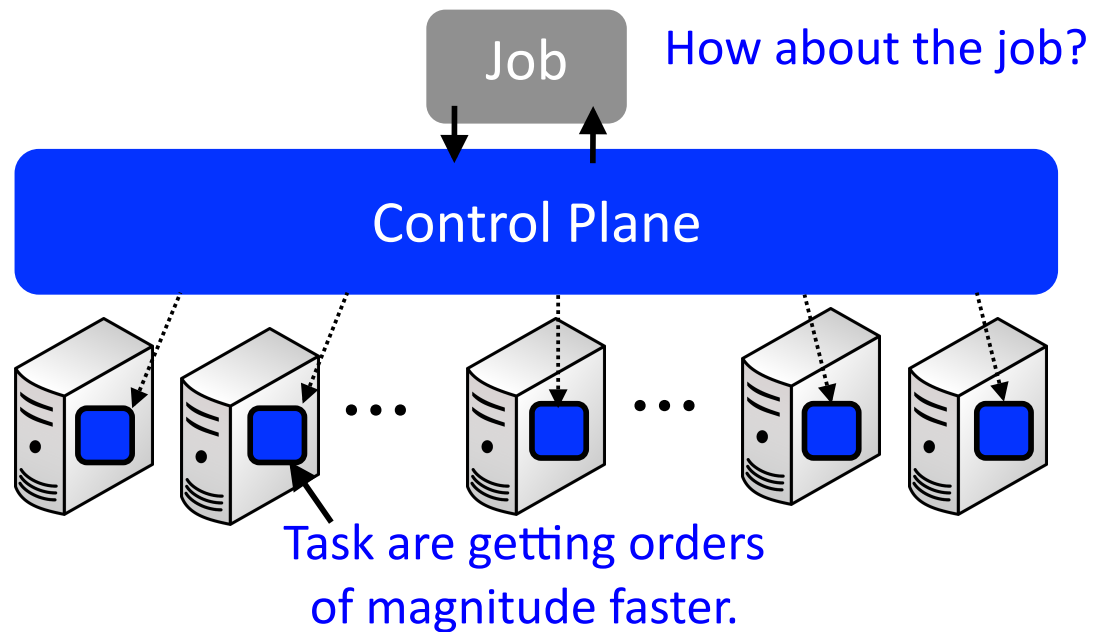


- One iteration of logistic regression over a data set of size 64MB.
- Tasks implemented efficiently, could run **50x** faster.

Individual tasks are getting faster.

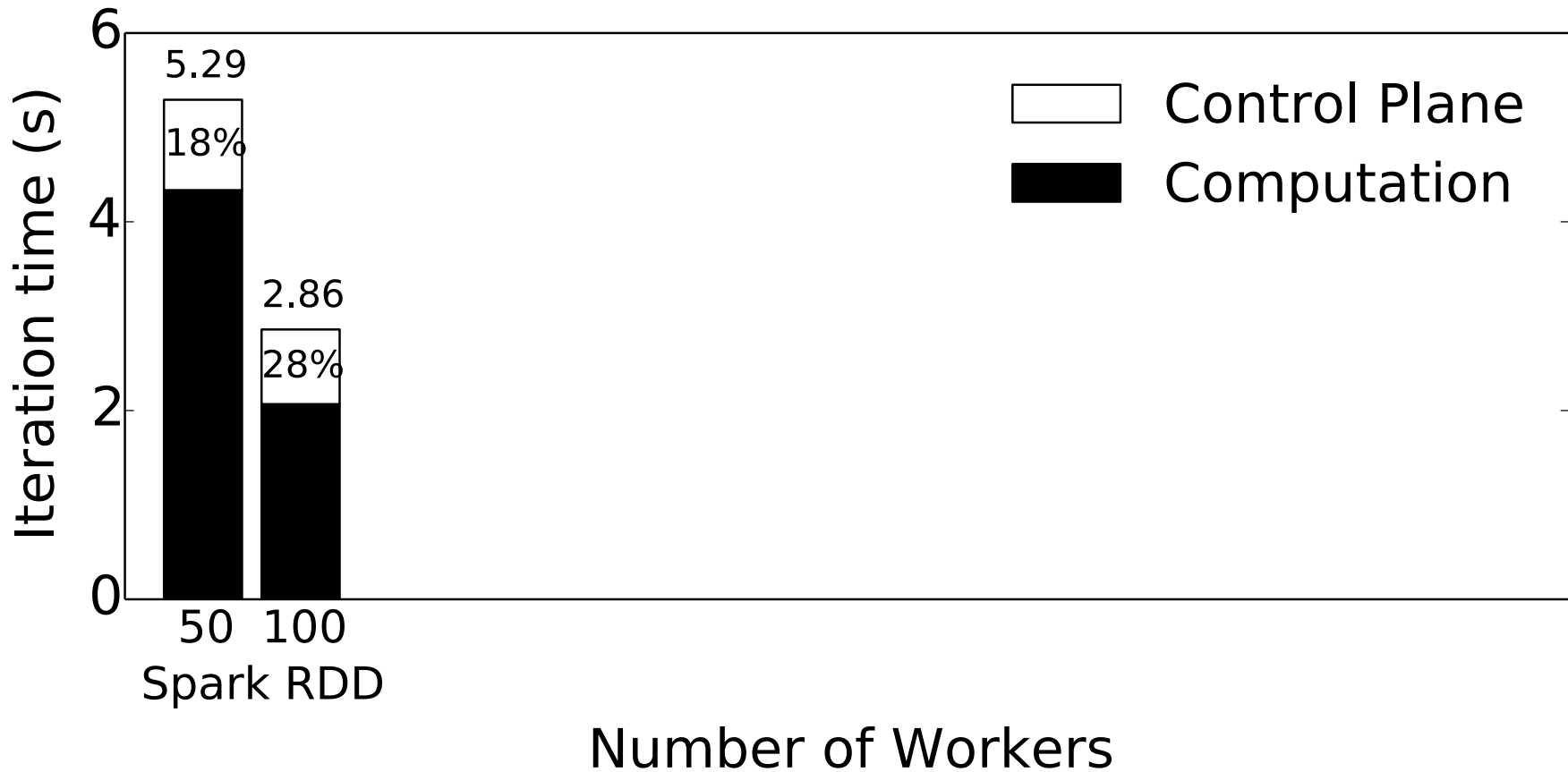
But does it necessarily mean that  
job completion time is getting shorter?

# Cloud Frameworks



# Control Plane

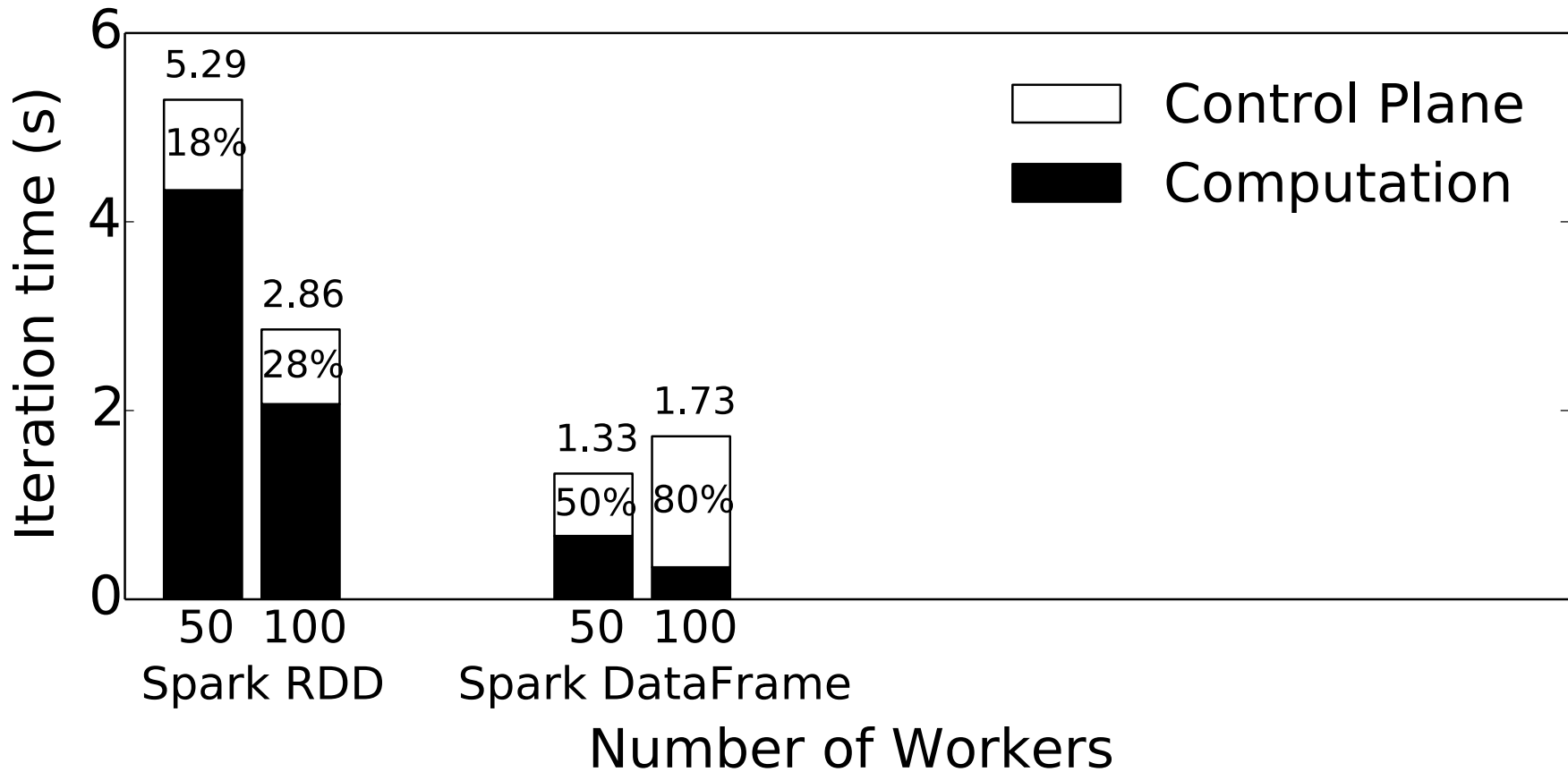
## The New Bottleneck



- Logistic regression over a data set of size 100GB.
- Classic Spark used to be **CPU-bound**.

# Control Plane

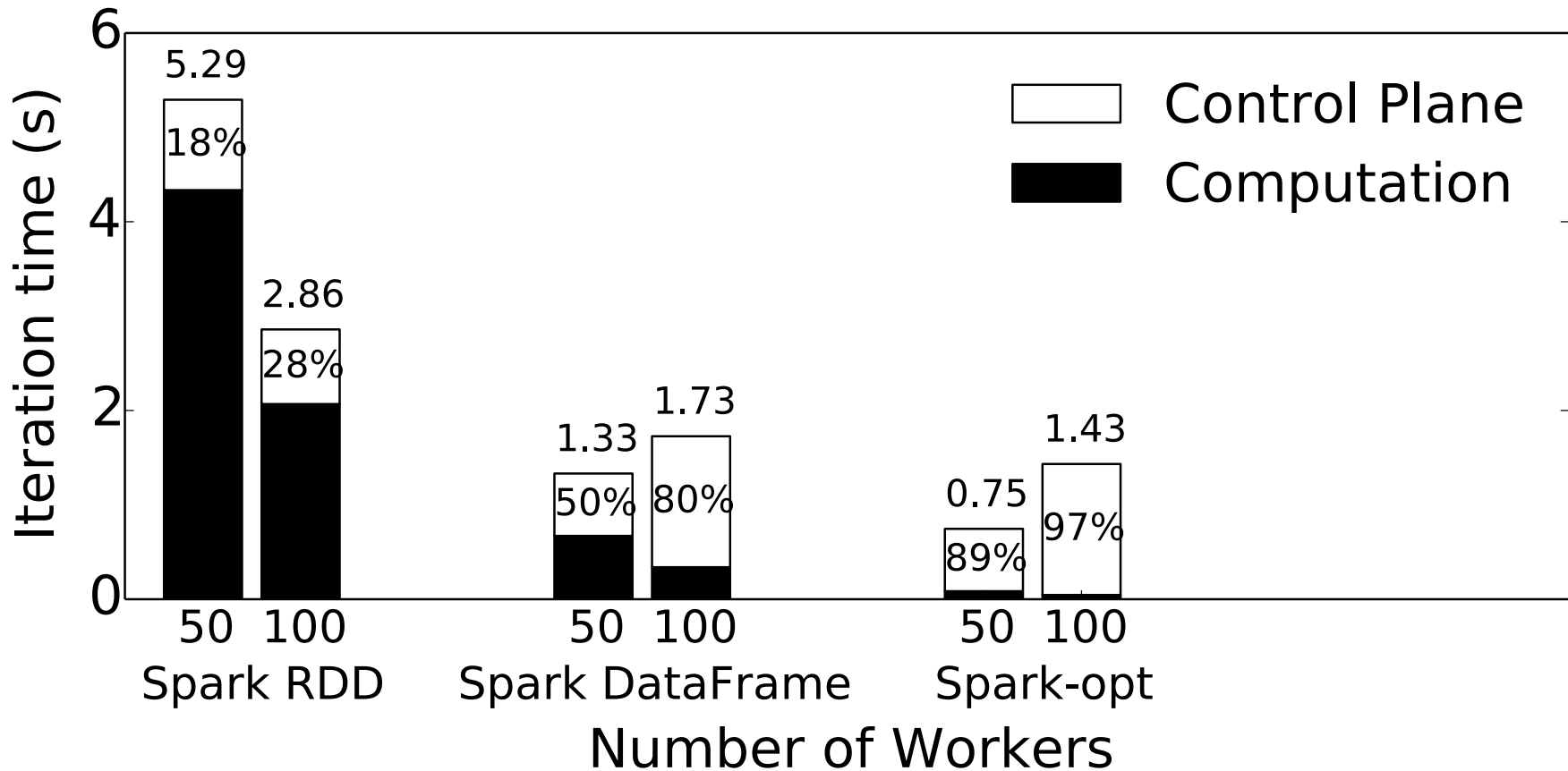
## The New Bottleneck



- Logistic regression over a data set of size 100GB.
- Spark 2.0 is already **control-bound**.

# Control Plane

## The New Bottleneck



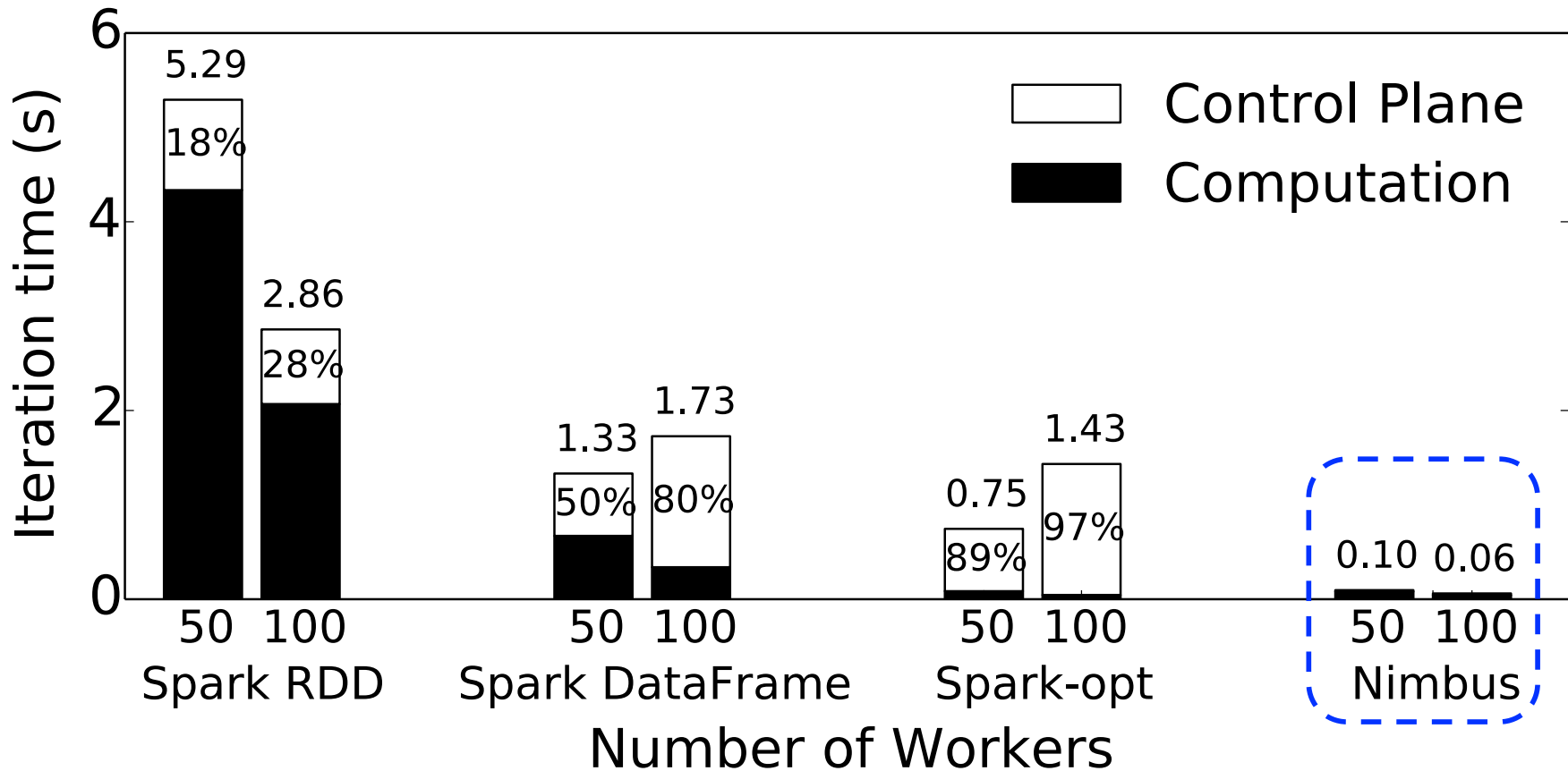
- Logistic regression over a data set of size 100GB.
- Spark-opt: hypothetical case where Spark runs tasks as fast as C++.

Control plane is the emerging bottleneck for the cloud computing frameworks.



# Control Plane

## The New Bottleneck



- Logistic regression over a data set of size 100GB.
- **Nimbus** with **execution templates** scales almost linearly.

# Contributions

- Demonstrating how the **control plane** is the emerging **bottleneck** for data analytics frameworks.
- **Execution Templates** as an abstraction for the control plane of cloud computing frameworks, that enables orders of magnitude higher task throughput, while keeping the fine-grained, flexible scheduling.
- The design, implementation, and evaluation of **Nimbus**, a distributed cloud computing framework that embeds execution templates.
- A demonstration of a single-core **graphical simulation** that Nimbus **automatically distributes** in the cloud showing execution templates in practice for complex applications.

# This talk

- Control Plane: the Emerging Bottleneck
- Design Scope of the Control Plane
- Execution Templates
- Nimbus: a Framework with Templates
- Evaluation

# This talk

- Control Plane: the Emerging Bottleneck
- Design Scope of the Control Plane
- Execution Templates
- Nimbus: a Framework with Templates
- Evaluation

# Cloud Frameworks Design

- Currently, there are two approaches:
  1. Centralized control model.
    - Controller generates and assigns tasks to the worker.
    - Limited task throughput, but reactive scheduling.
  2. Distributed data flow model.
    - Nodes generate and spawn tasks locally.
    - Great scalability, but static scheduling.

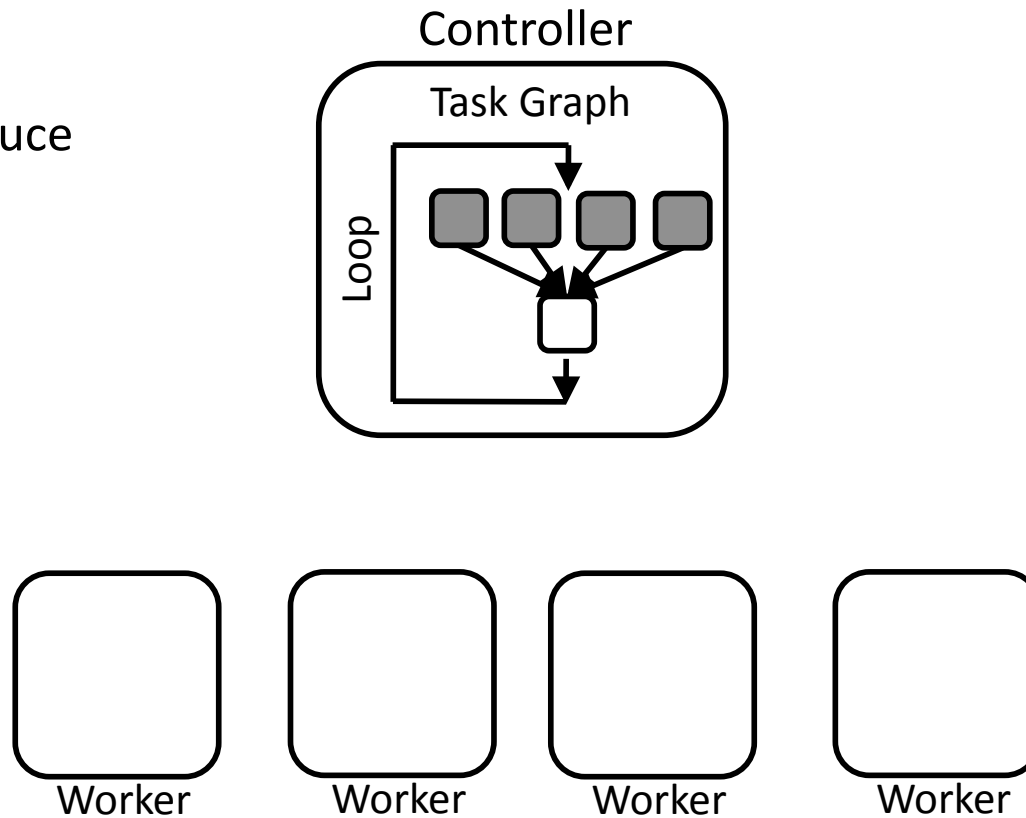


## Design Spectrum

Centralized Controllers

Distributed Controllers

- MapReduce
- Hadoop
- Spark



- Controller centrally schedules and spawns tasks.





# Distributed Controllers







## A long horizontal double-headed arrow, consisting of a thick black line with arrowheads at both ends, spanning most of the width of the page.

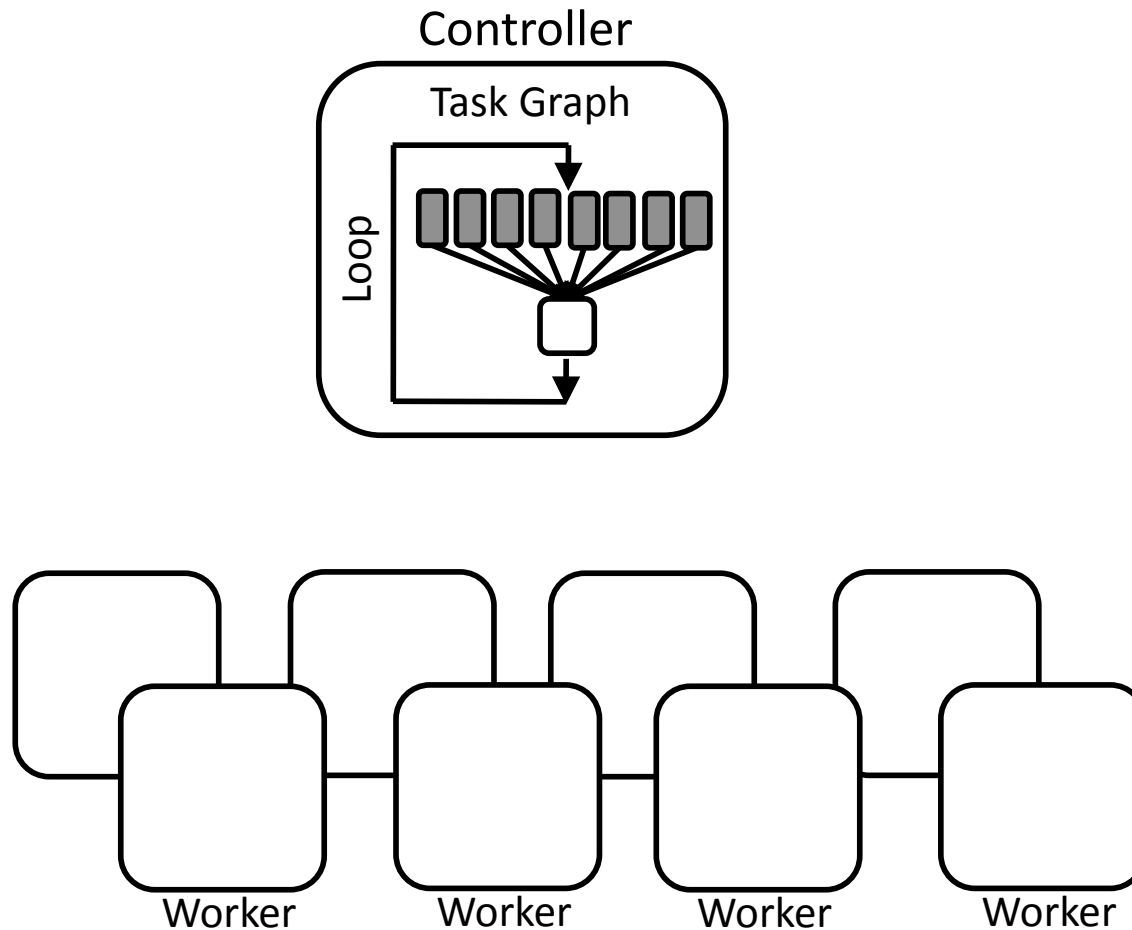
## Distributed Controllers



## Design Spectrum

Centralized Controllers

Distributed Controllers



- But controller bottlenecks at scale.

# Distributed Controllers

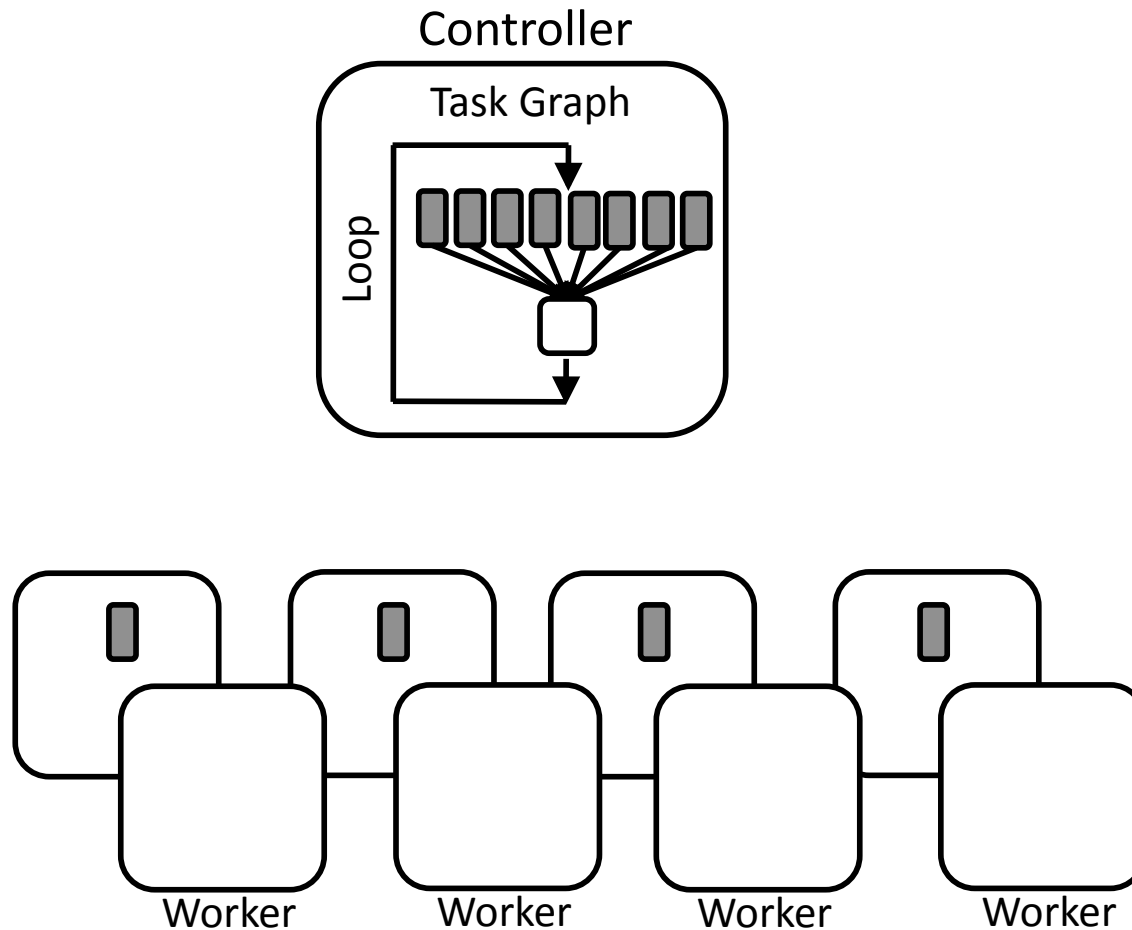




## Design Spectrum

Centralized Controllers

Distributed Controllers



- But controller bottlenecks at scale.

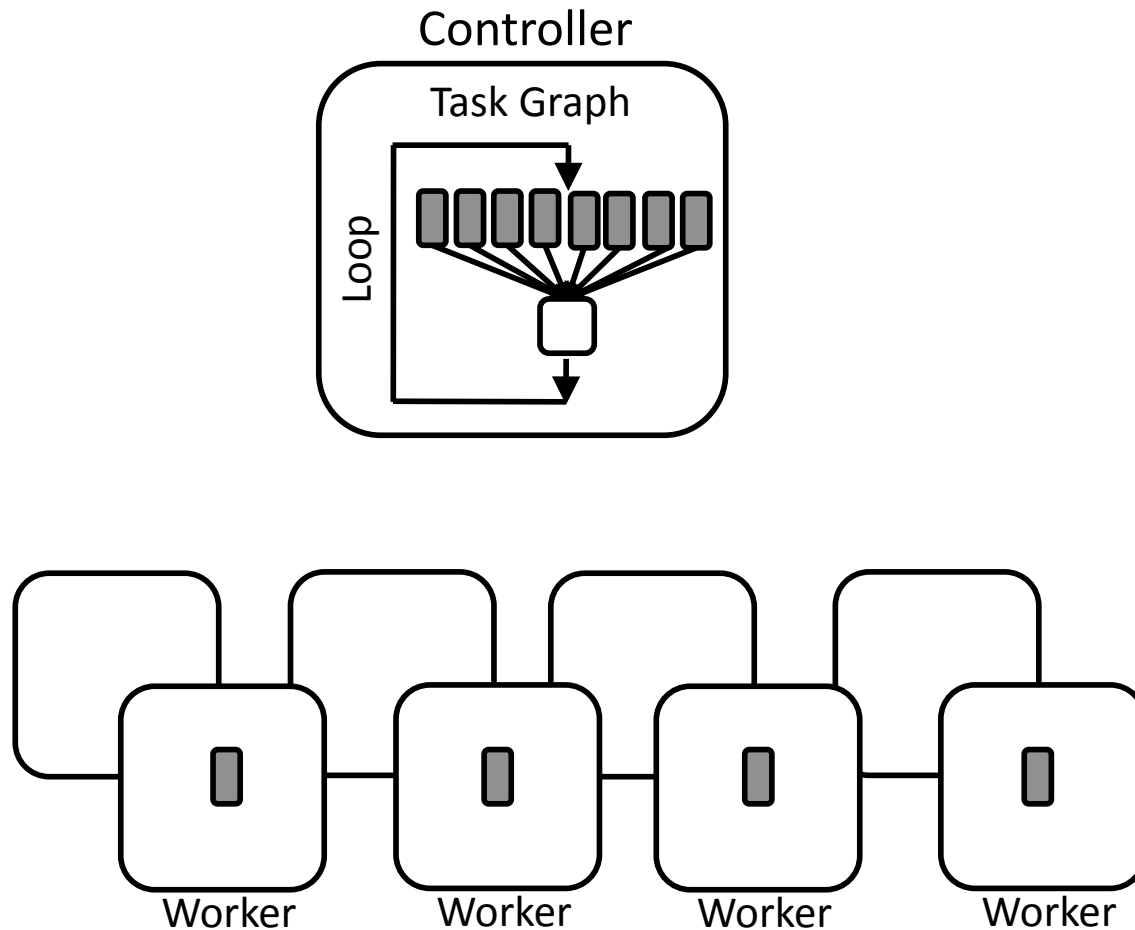
## Distributed Controllers



## Design Spectrum

Centralized Controllers

Distributed Controllers



- But controller bottlenecks at scale.

## Design Spectrum

Centralized Controllers

Distributed Controllers



- Logistic regression over a data set of size 100GB in Spark 2.0 MLlib.
- Control Plane **bottlenecks at scale**, generating and spawning tasks.

## A horizontal line with arrowheads at both ends, pointing left and right.

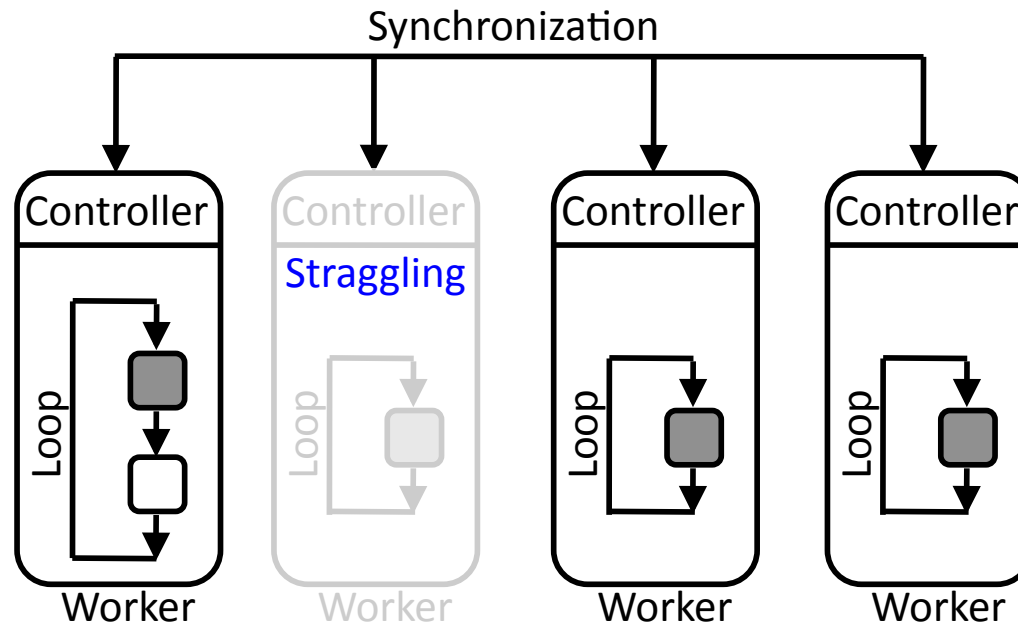
# Distributed Controllers



- 34

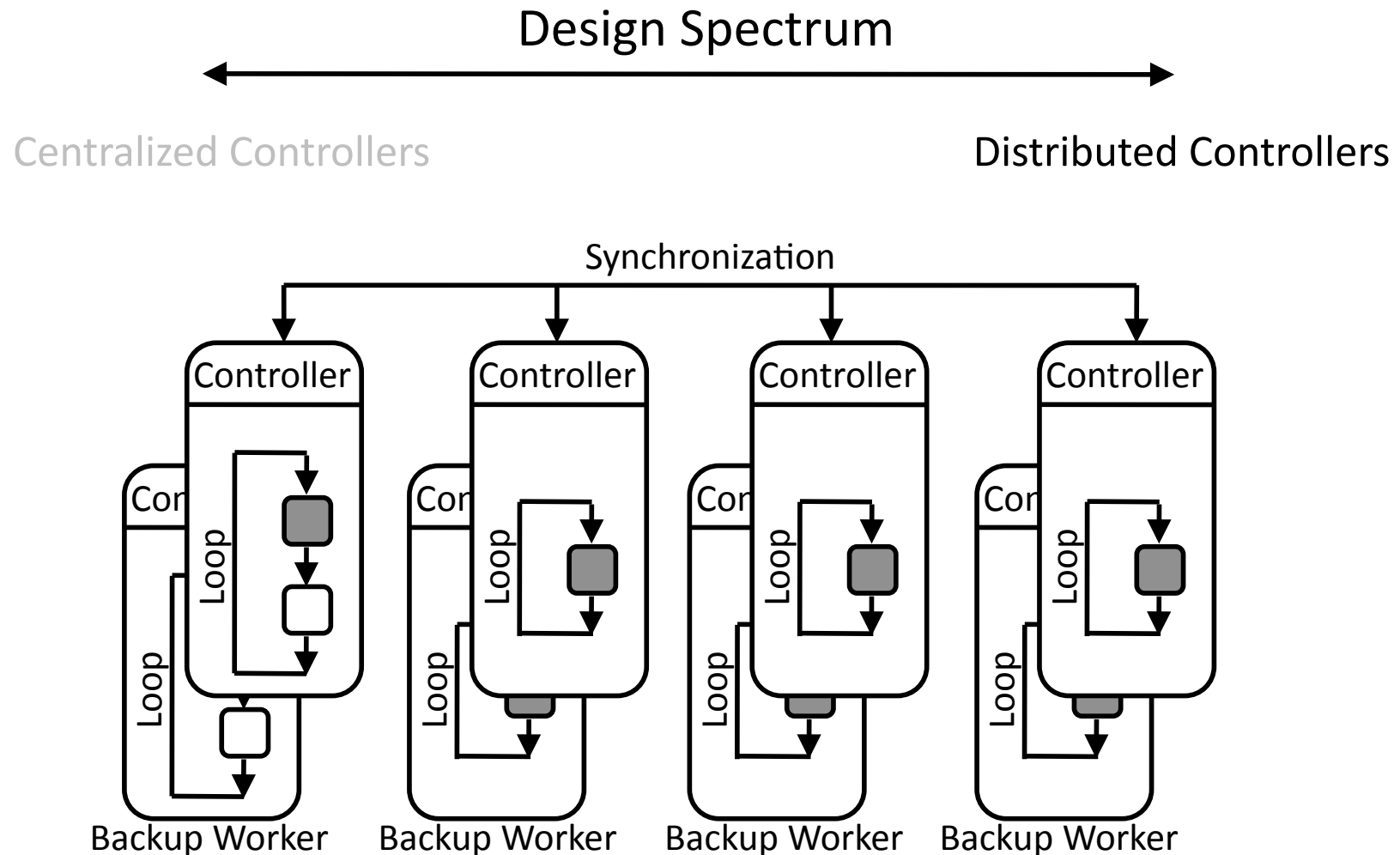


# Distributed Controllers



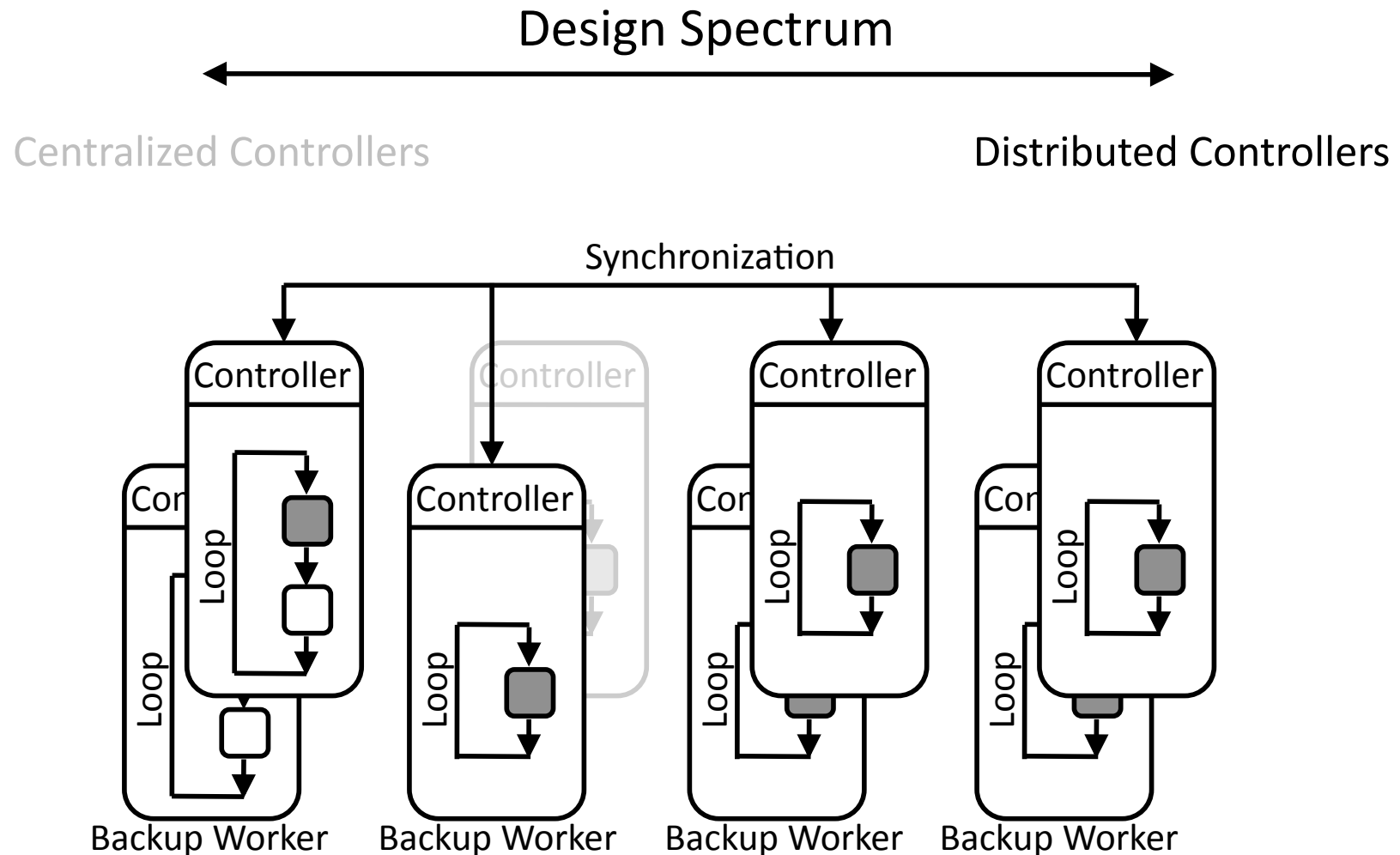
- But, the scheduling is static.
  - The progress speed is bound to the speed of the slowest node.
  - Any change requires stopping all nodes and installing new data flow.





- In practice the straggler mitigation is only **proactive**:
  - Avoiding stragglers by meticulous engineering work.
  - Launching backup workers (at least doubling the resources).





- In practice the straggler mitigation is only **proactive**:
  - Avoiding stragglers by meticulous engineering work.
  - Launching backup workers (at least doubling the resources).

# Design Space

## Summary

Control Plane Design	Example Framework	Task Throughput	Task Scheduling
Centralized	MapReduce	Low	Dynamic
	Hadoop		
	Spark		
Distributed	Naiad	High	Static
	TensorFlow		

# Design Space

## Summary

Control Plane Design	Example Framework	Task Throughput	Task Scheduling
Centralized	MapReduce	Low	Dynamic
	Hadoop		
	Spark		
Distributed	Naiad	High	Static
	TensorFlow		

- We would like to have the best of both worlds:
- High task throughput for fast computations.
  - Dynamic, fine-grained scheduling decisions.

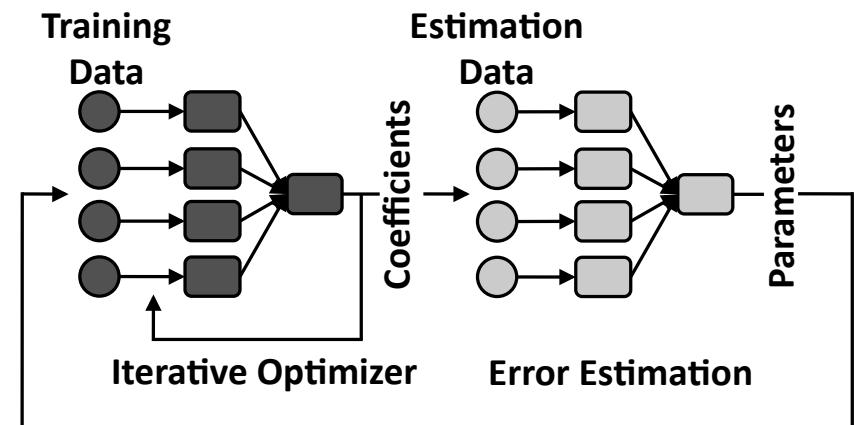
# Repetitive Patterns

- Advanced data analytics are iterative in nature.
  - Machine learning, graph processing, image recognition, etc.
- This results in repetitive patterns in the control plane.
  - Similar tasks execute with minor differences.

# Repetitive Patterns

- Advanced data analytics are iterative in nature.
  - Machine learning, graph processing, image recognition, etc.
- This results in repetitive patterns in the control plane.
  - Similar tasks execute with minor differences.

```
while (error > threshold_e) {  
  while (gradient > threshold_g) {  
    // Optimization code block  
    gradient = Gradient(tdata, coeff, param)  
    coeff += gradient  
  }  
  // Estimation code block  
  error = Estimate(edata, coeff, param)  
  param = update_model(param, error)  
}
```



# This talk

- Control Plane: the Emerging Bottleneck
- Design Scope of the Control Plane
- Execution Templates
- Nimbus: a Framework with Templates
- Evaluation



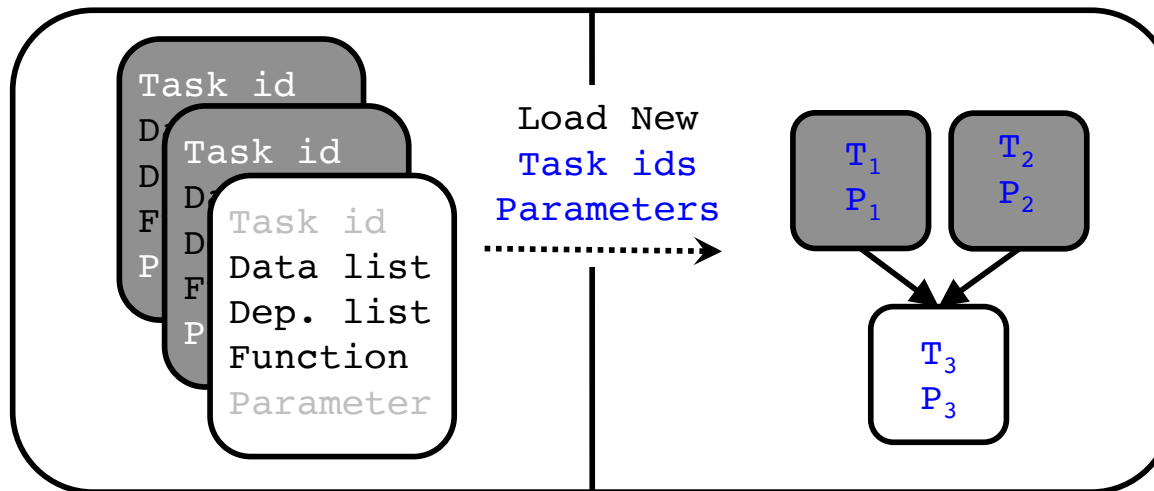
# Execution Templates

- Tasks are cached as **parameterizable blocks** on nodes.
- Instead of assigning the tasks from scratch, templates are **instantiated** by filling in only changing parameters.



# Execution Templates

- Tasks are cached as **parameterizable blocks** on nodes.
- Instead of assigning the tasks from scratch, templates are **instantiated** by filling in only changing parameters.

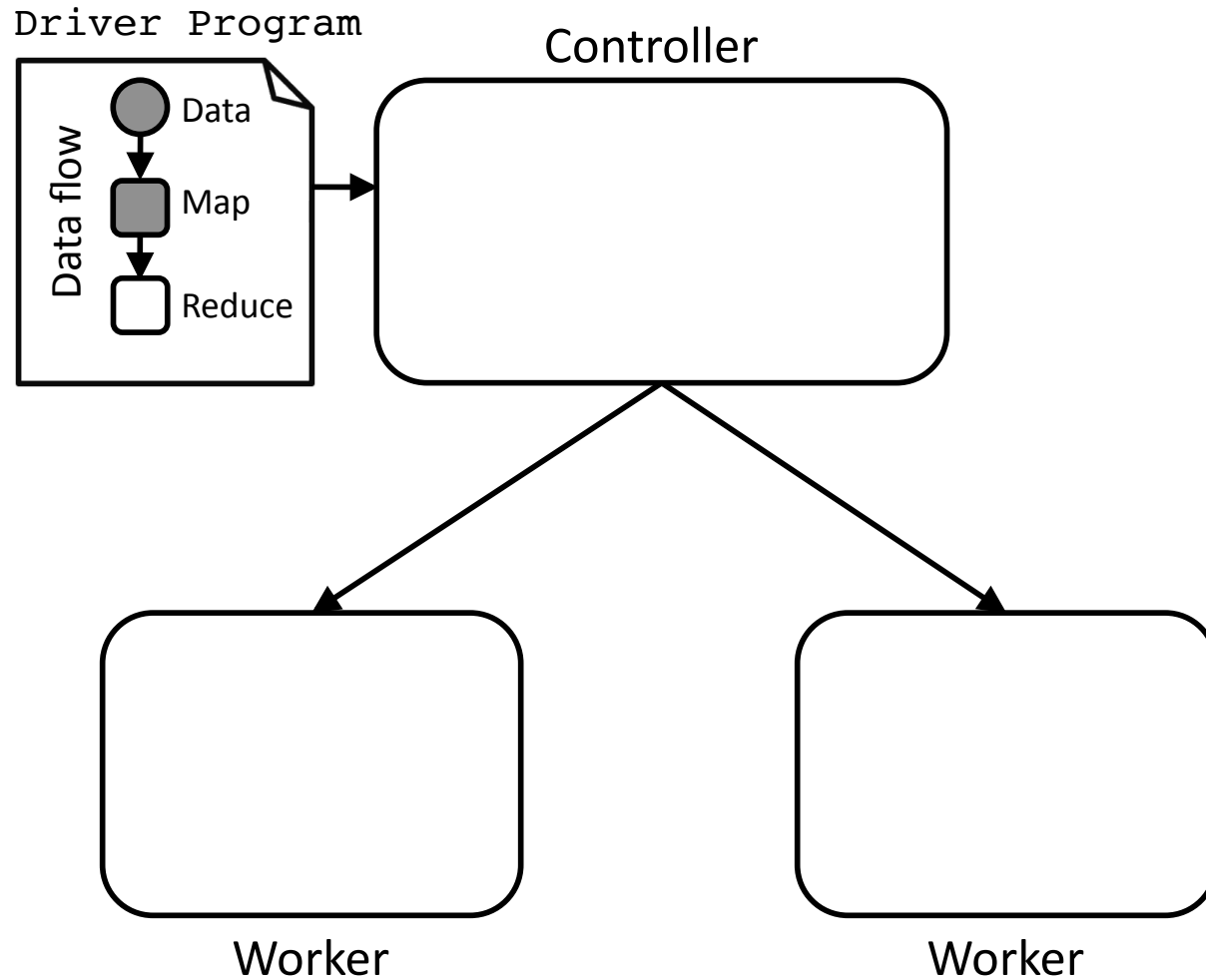


# Execution Templates

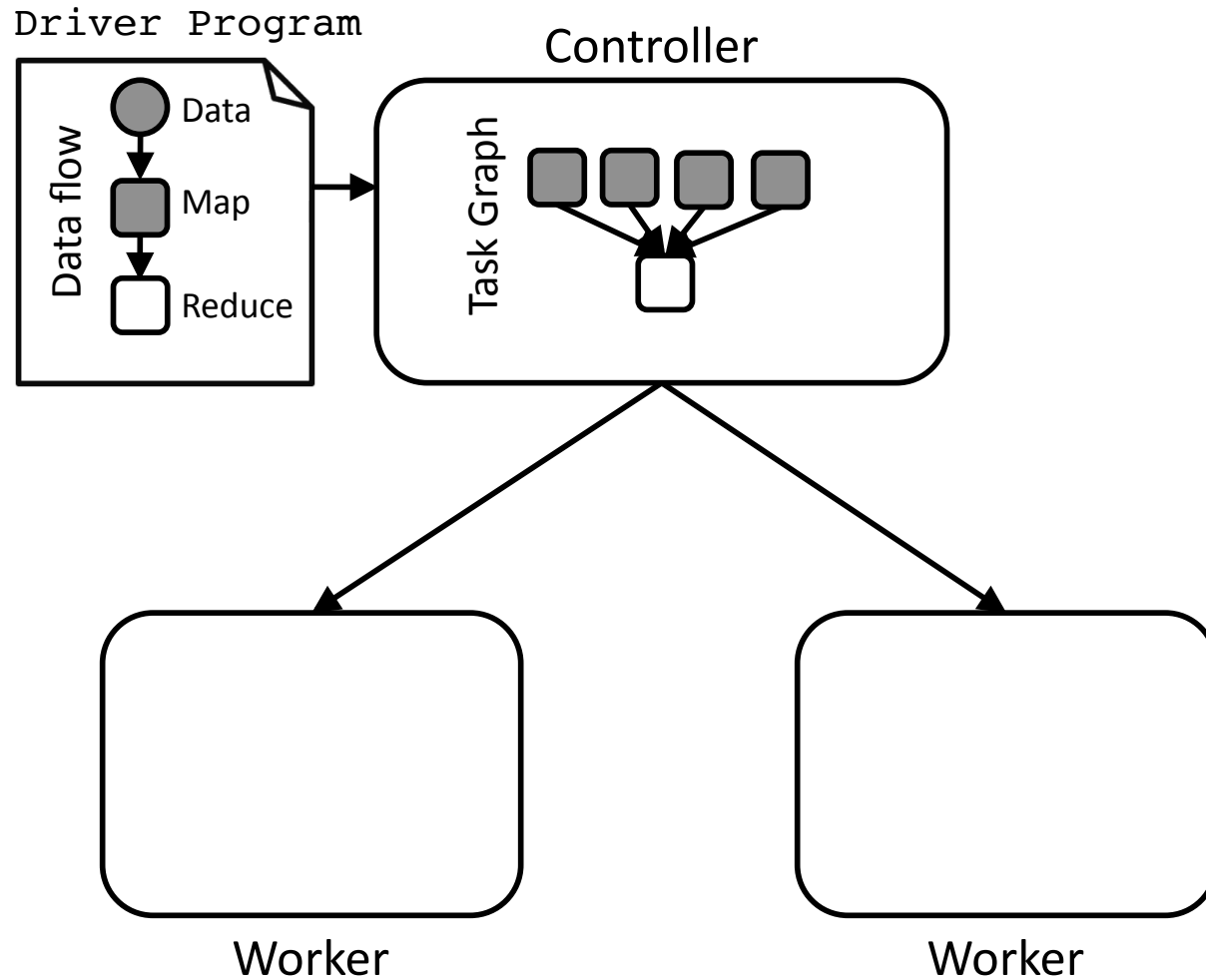
## Mechanisms Summary

- **Instantiation:** spawn a block of tasks without processing each task individually from scratch. It helps increase the **task throughput**.
- **Edits:** modifies the content of each template at the granularity of tasks. It enables fine-grained, **dynamic scheduling**.
- **Patches:** In case the state of the worker does not match the preconditions of the template. It enables **dynamic control flow**.

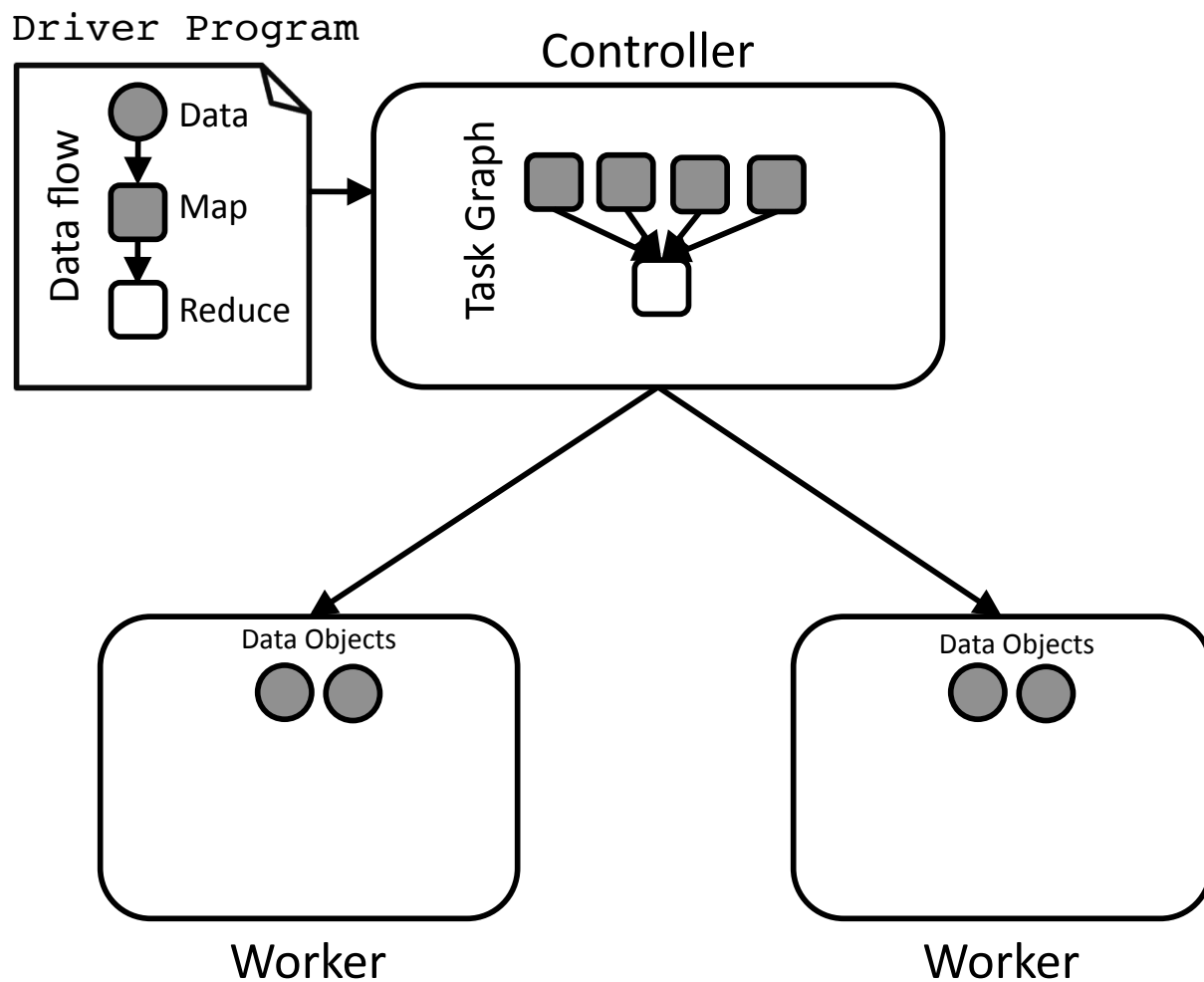
# Execution Model



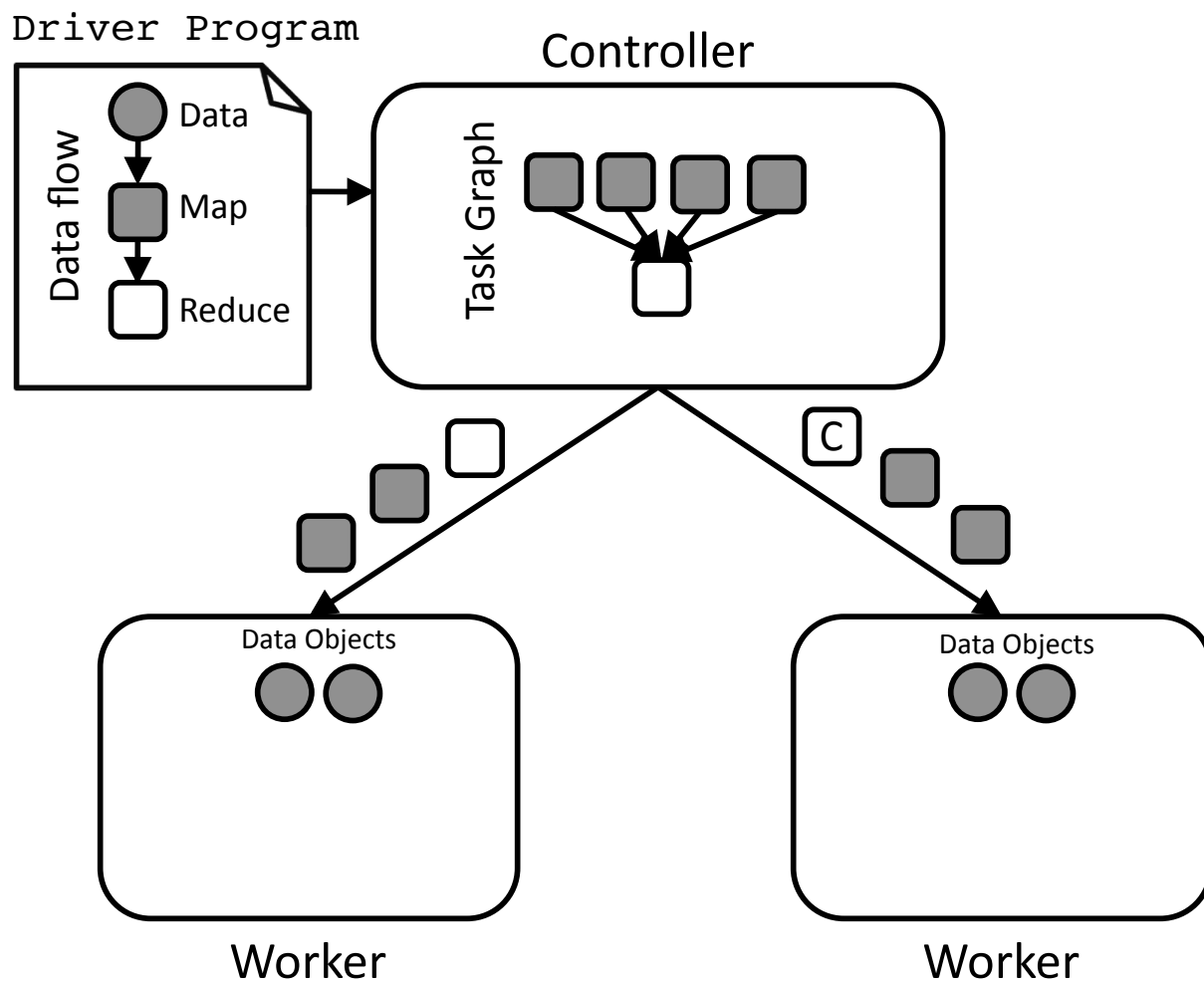
# Execution Model



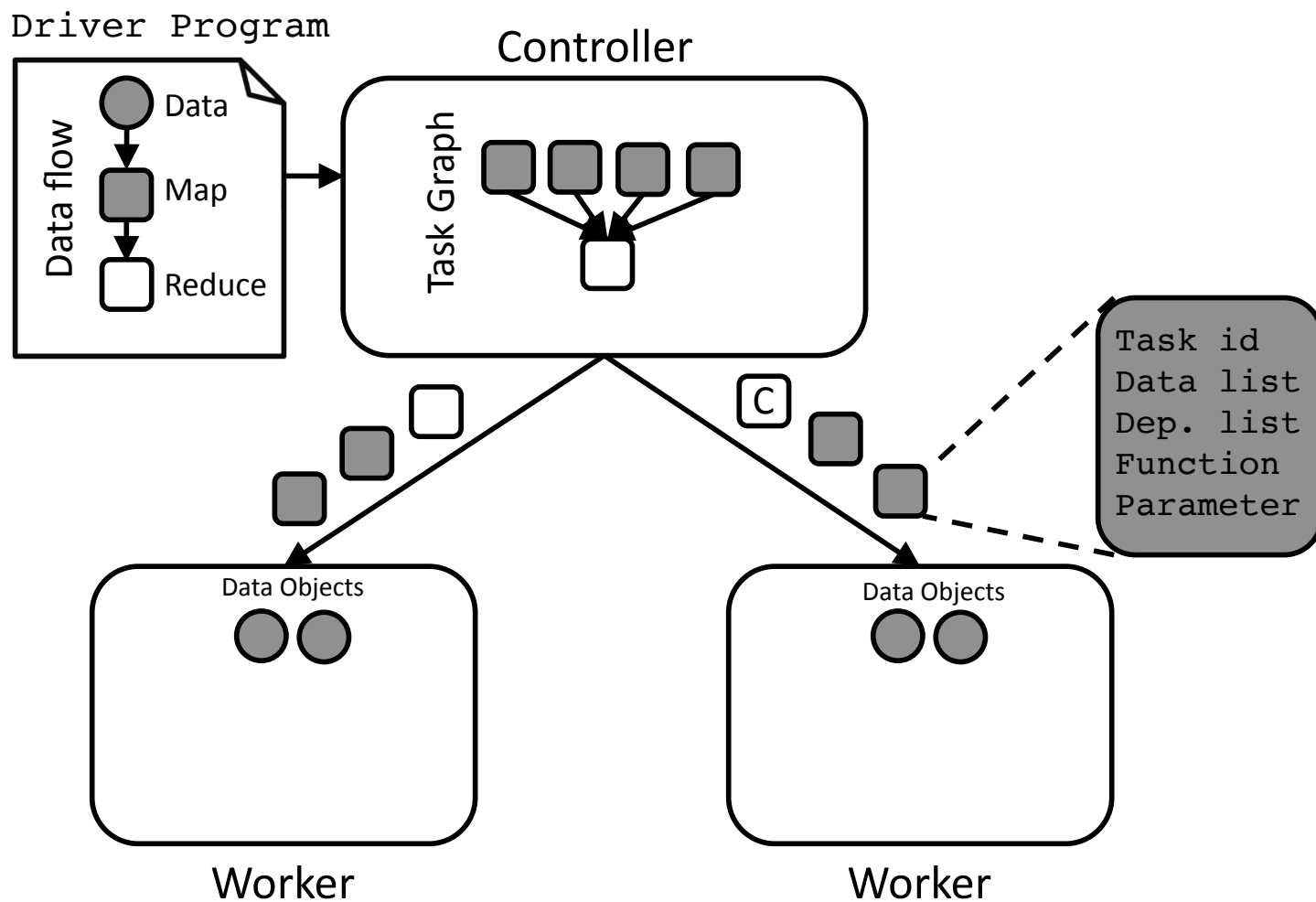
# Execution Model



# Execution Model

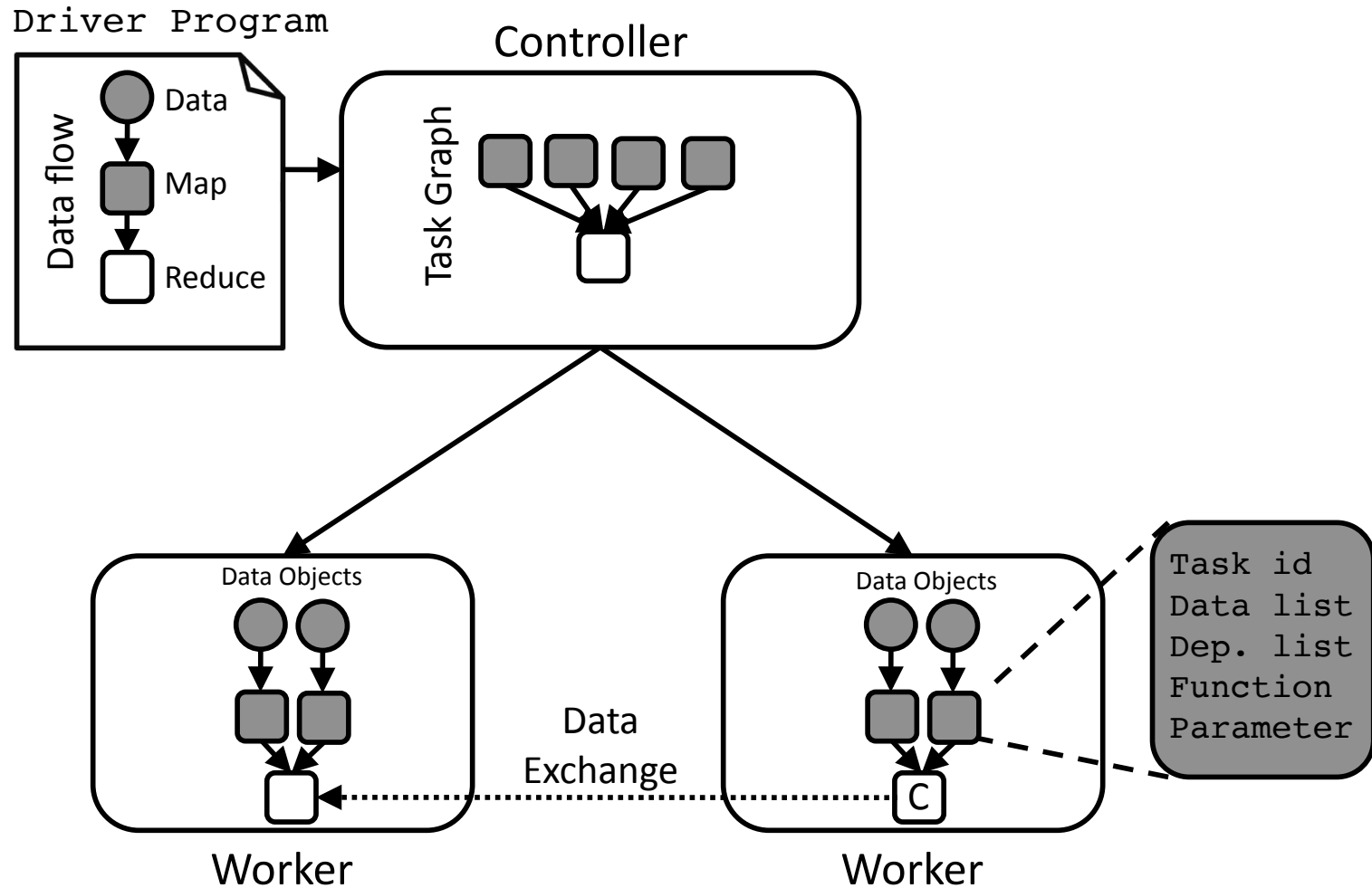


# Execution Model

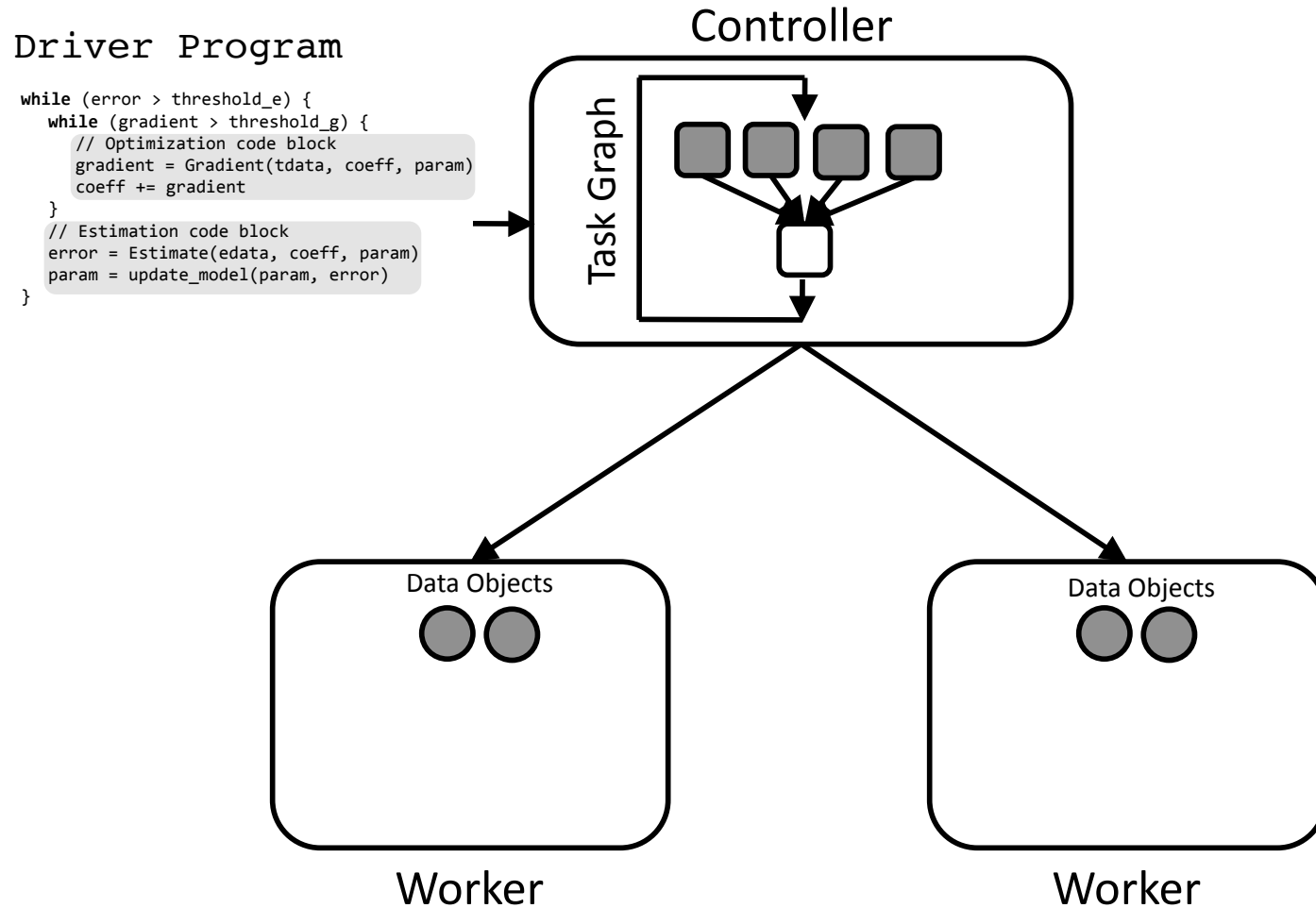




# Execution Model



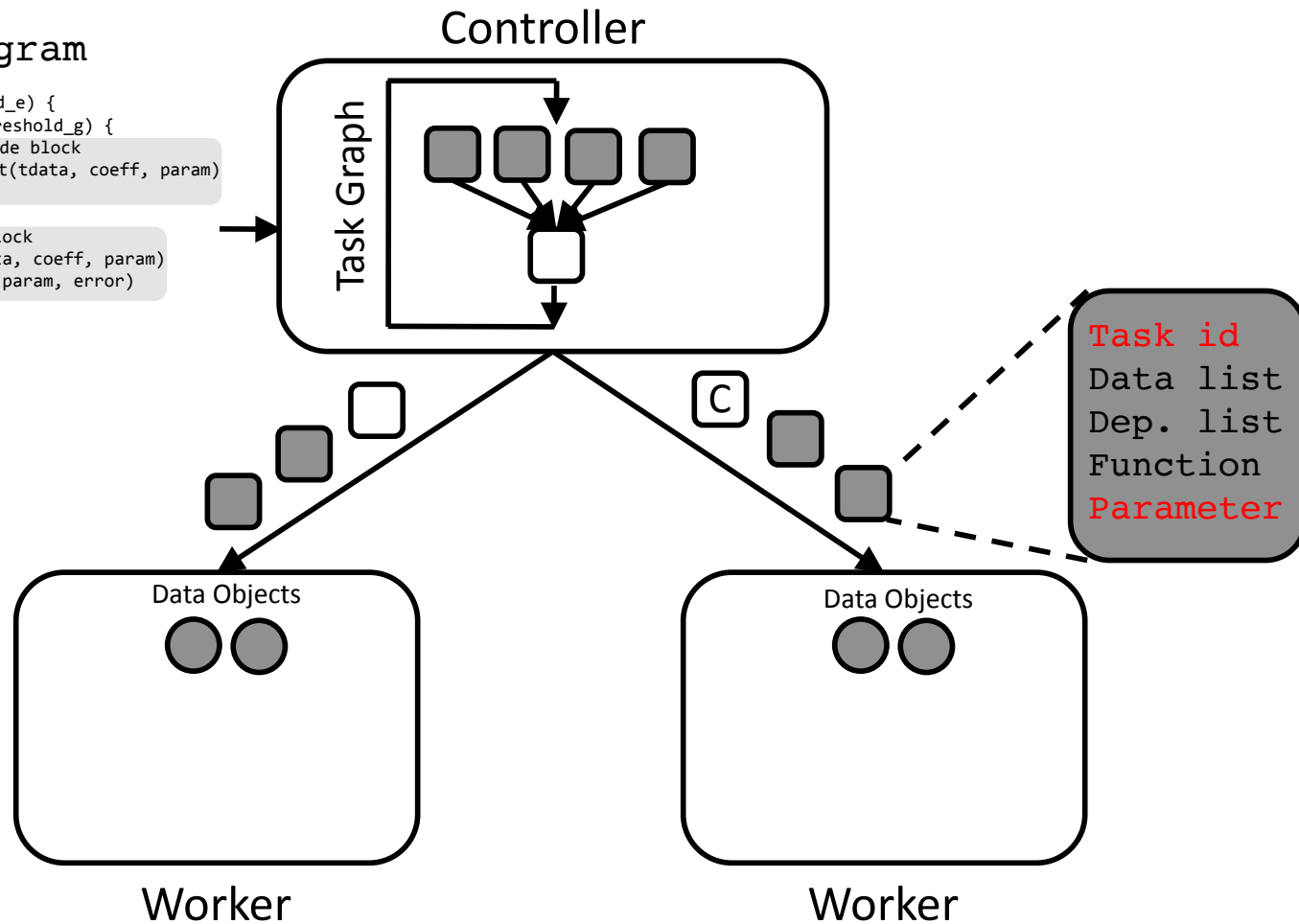
# Repetitive Patterns



# Repetitive Patterns

## Driver Program

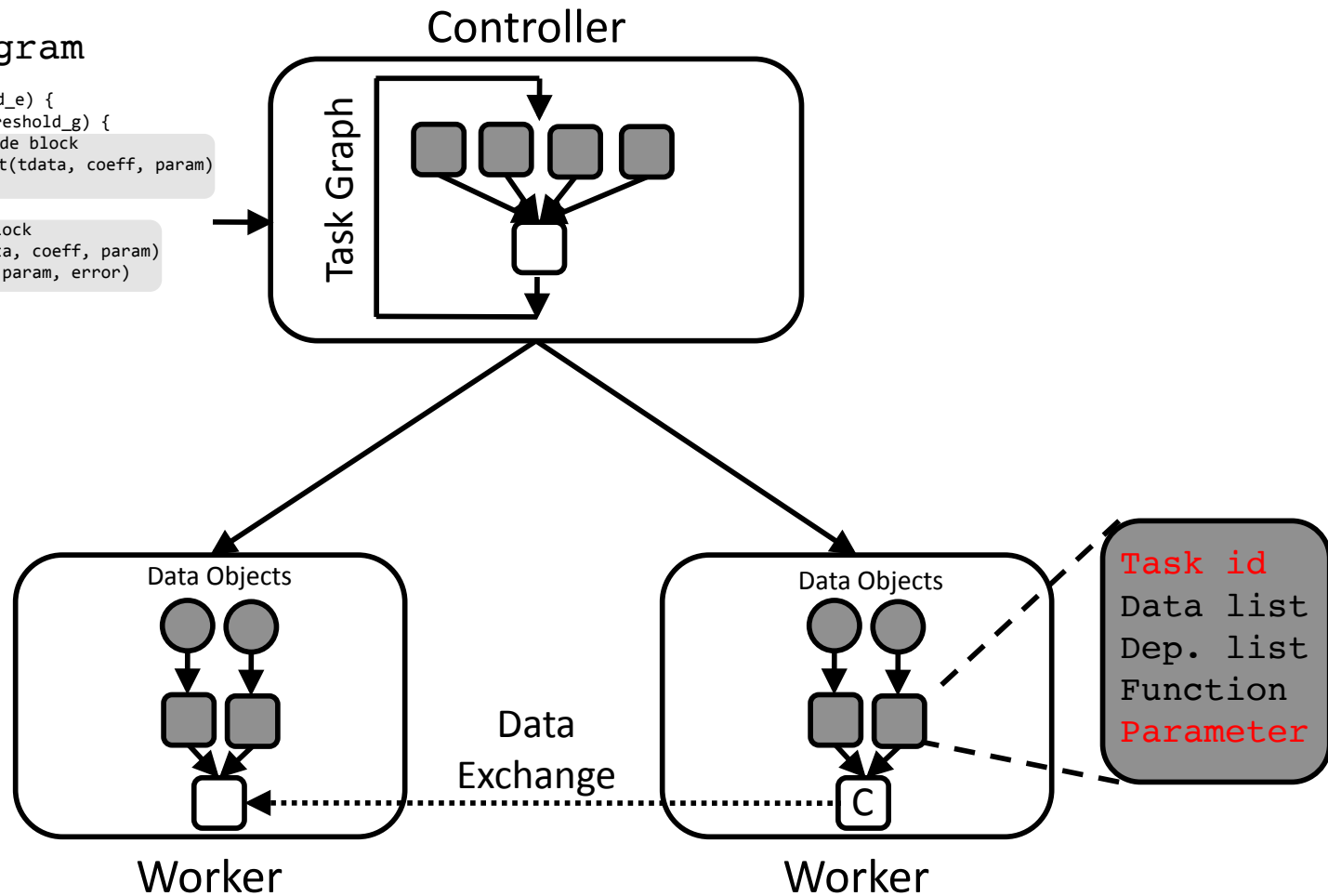
```
while (error > threshold_e) {  
  while (gradient > threshold_g) {  
    // Optimization code block  
    gradient = Gradient(tdata, coeff, param)  
    coeff += gradient  
  }  
  // Estimation code block  
  error = Estimate(edata, coeff, param)  
  param = update_model(param, error)  
}
```



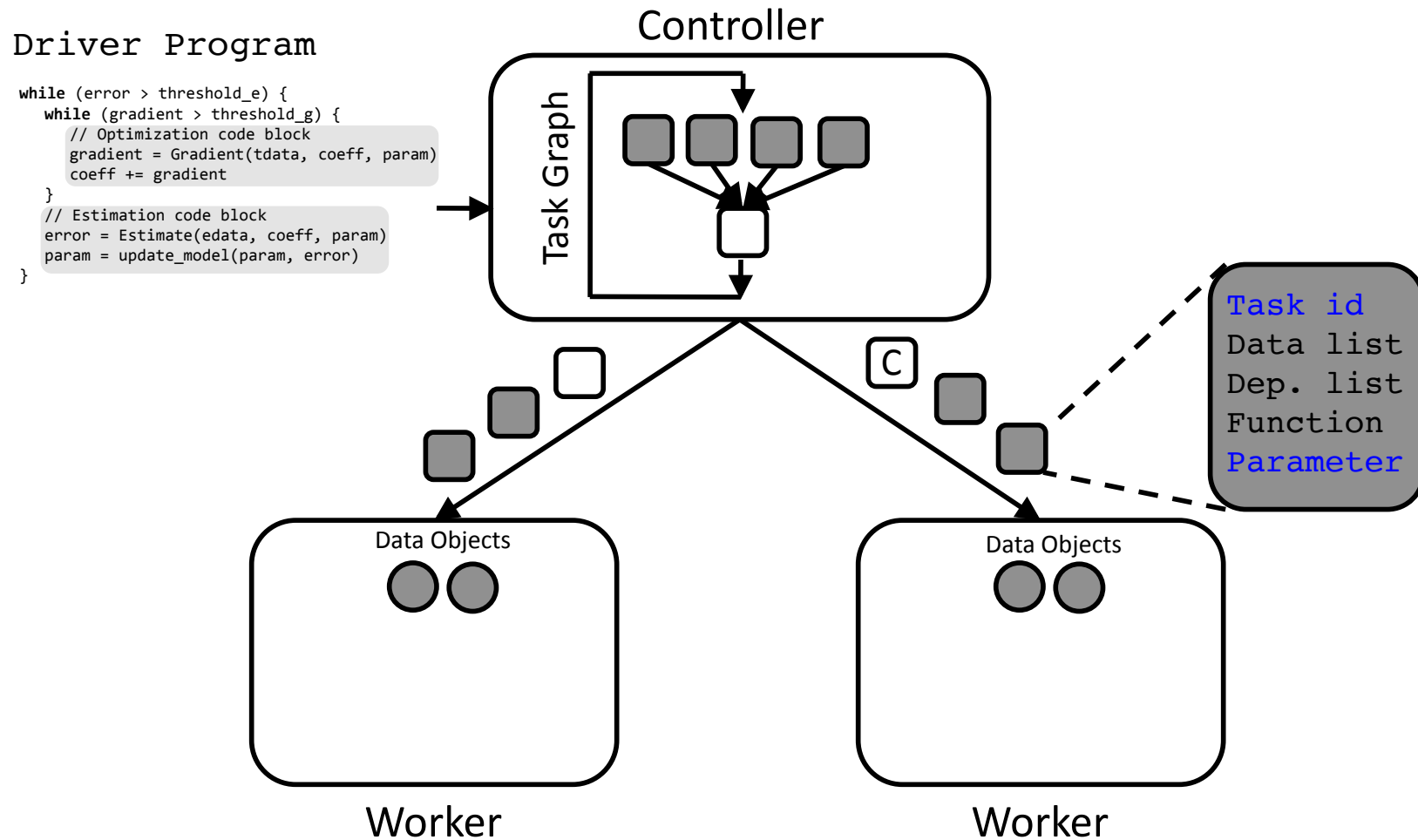
# Repetitive Patterns

## Driver Program

```
while (error > threshold_e) {  
  while (gradient > threshold_g) {  
    // Optimization code block  
    gradient = Gradient(tdata, coeff, param)  
    coeff += gradient  
  }  
  // Estimation code block  
  error = Estimate(edata, coeff, param)  
  param = update_model(param, error)  
}
```



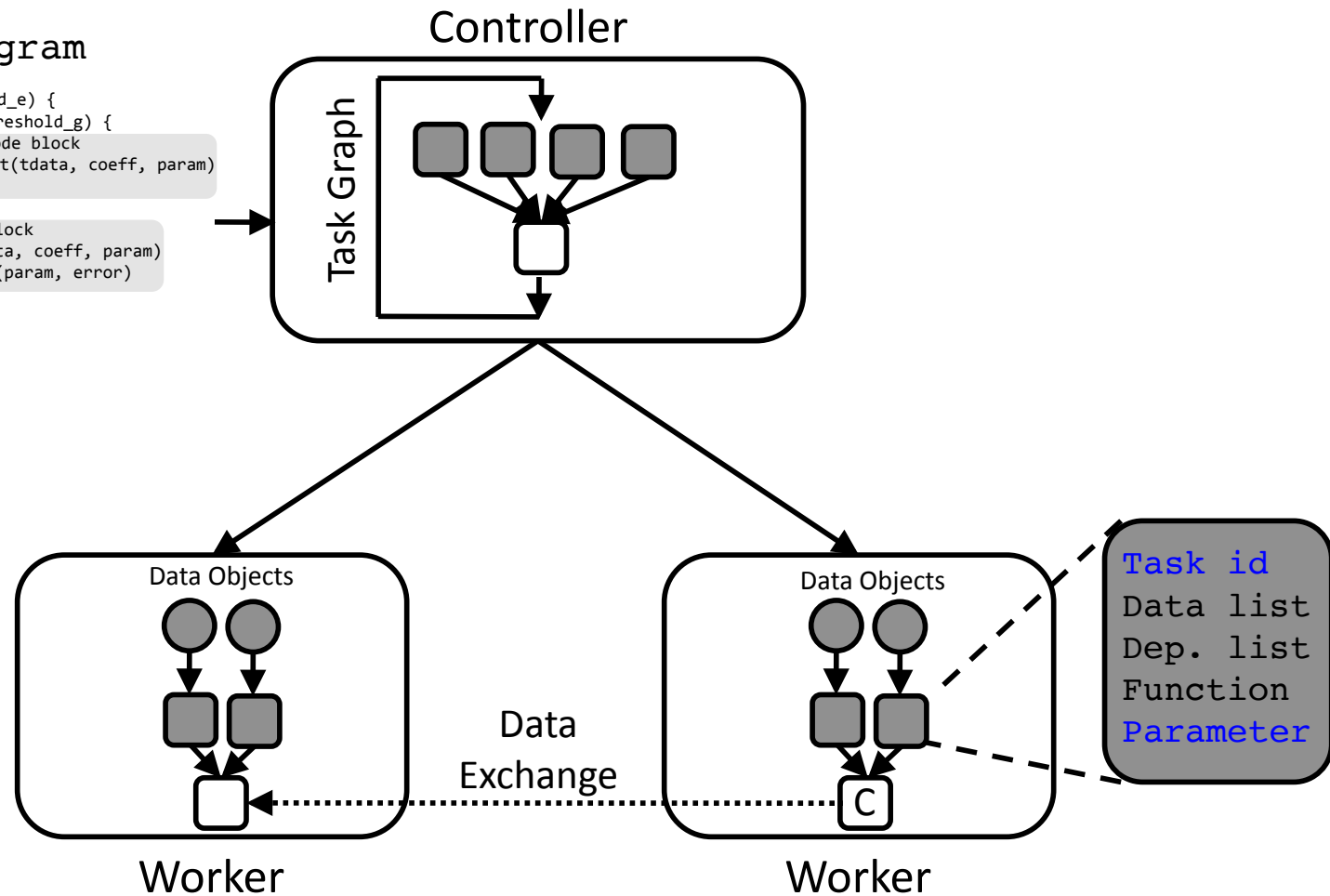
# Repetitive Patterns



# Repetitive Patterns

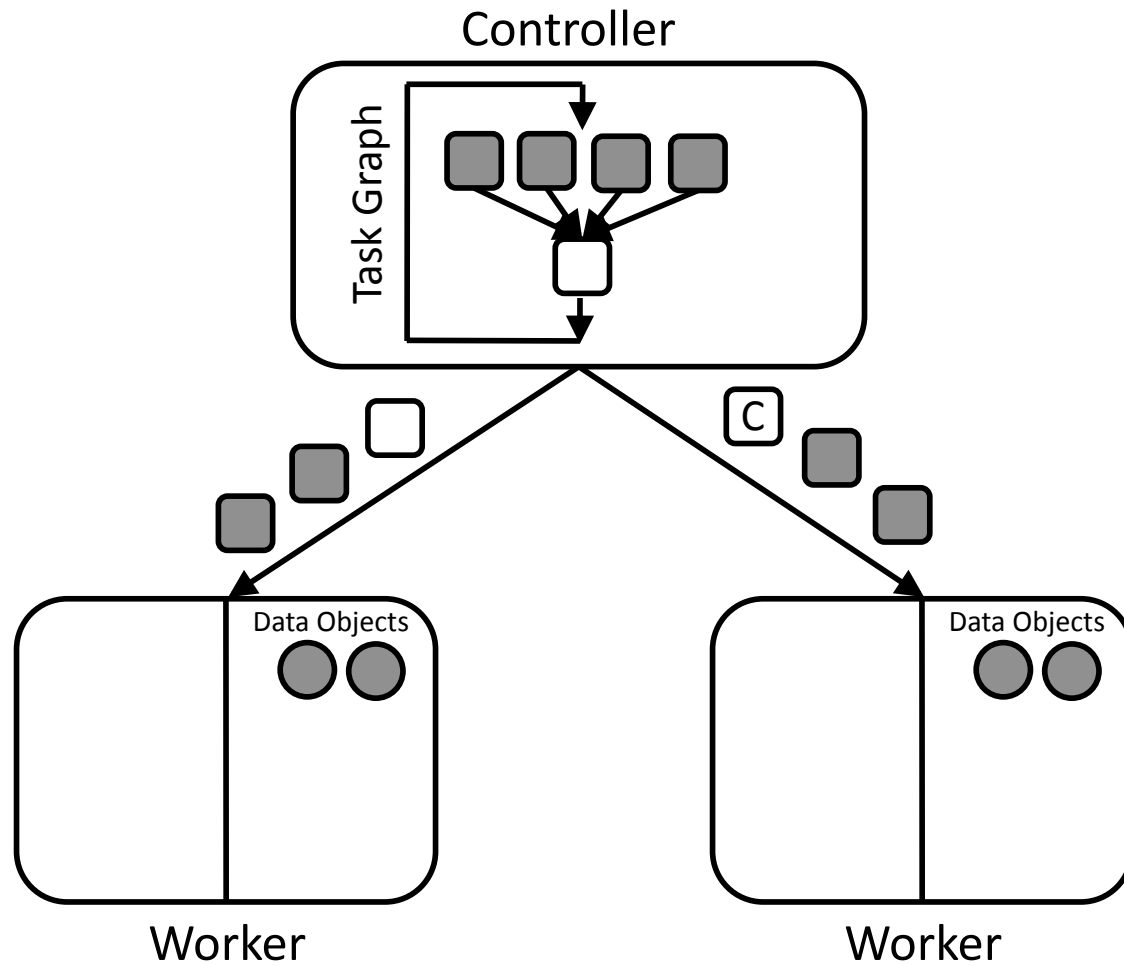
## Driver Program

```
while (error > threshold_e) {  
  while (gradient > threshold_g) {  
    // Optimization code block  
    gradient = Gradient(tdata, coeff, param)  
    coeff += gradient  
  }  
  // Estimation code block  
  error = Estimate(edata, coeff, param)  
  param = update_model(param, error)  
}
```



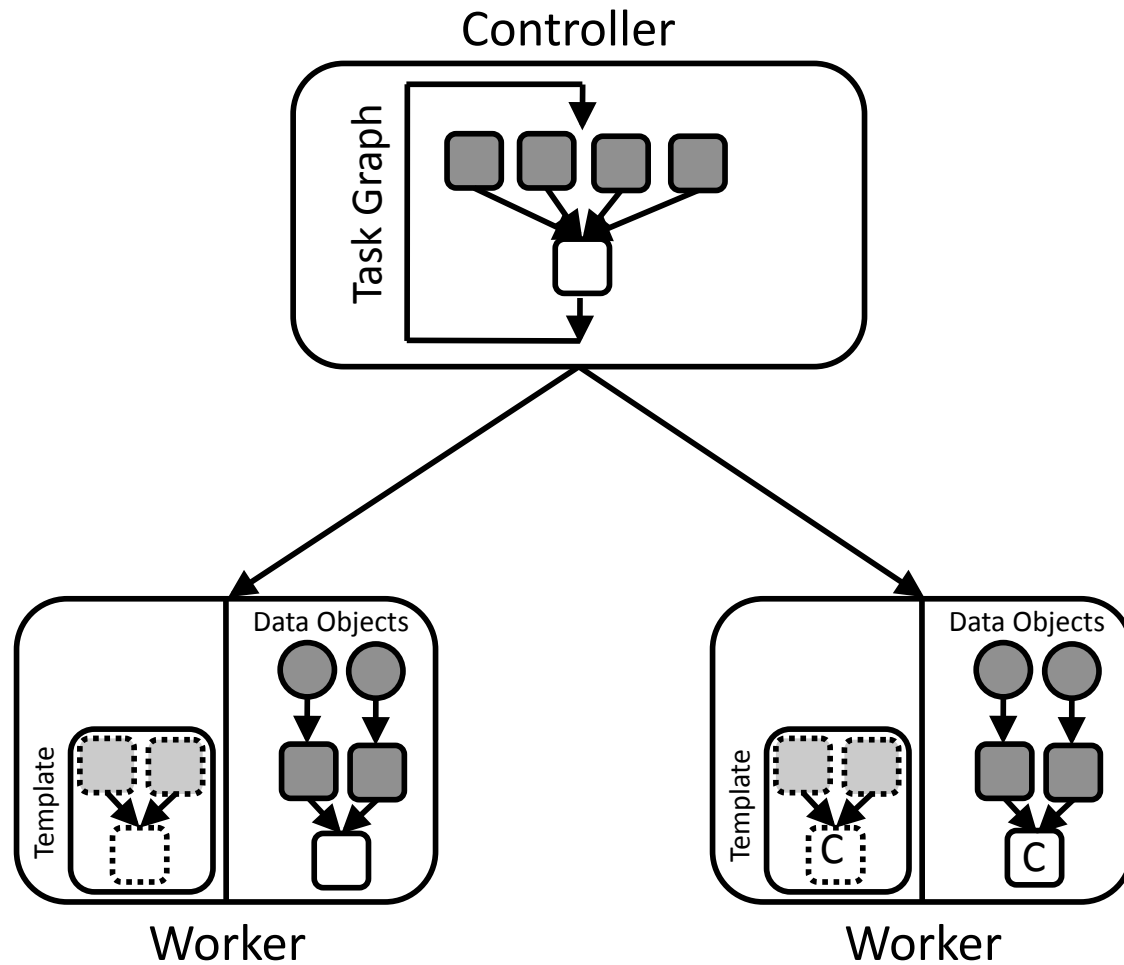
# Execution Templates

## Abstraction



# Execution Templates

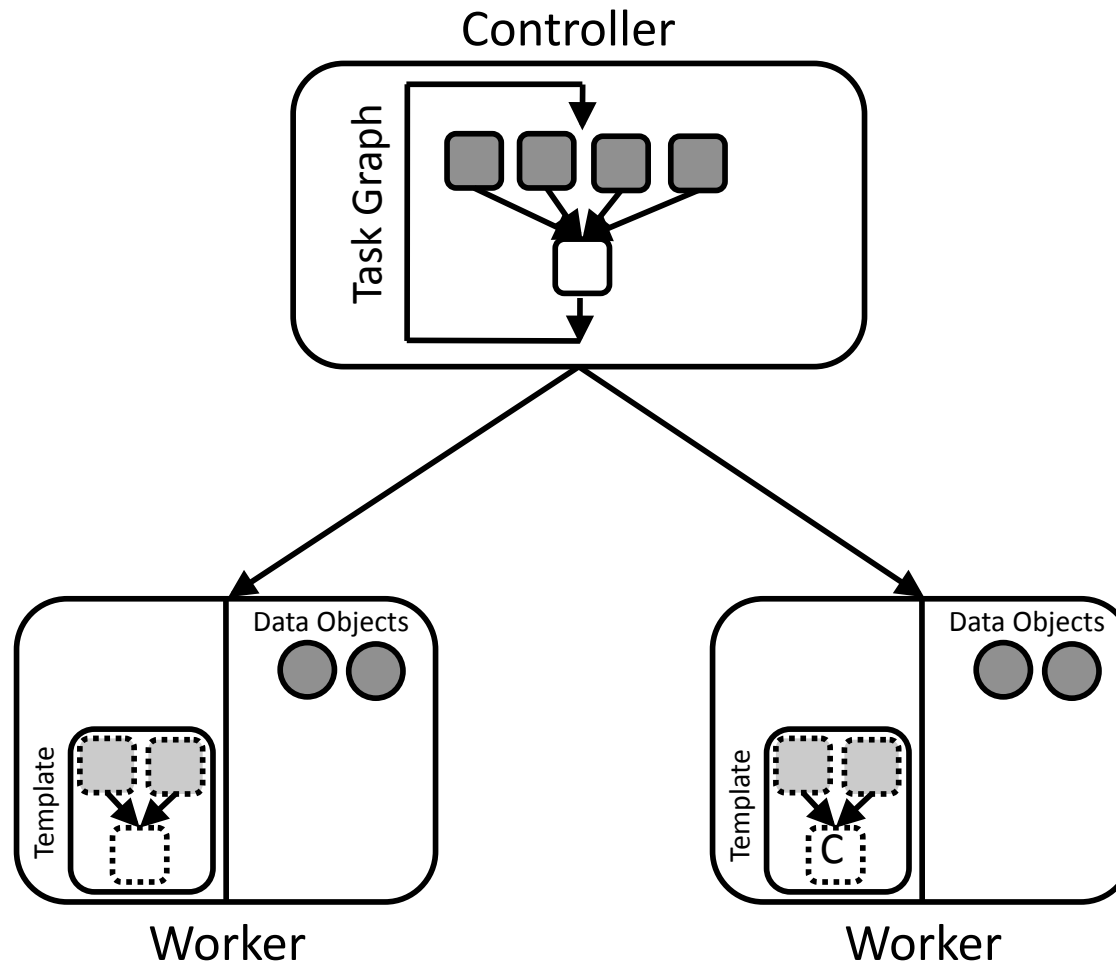
## Abstraction





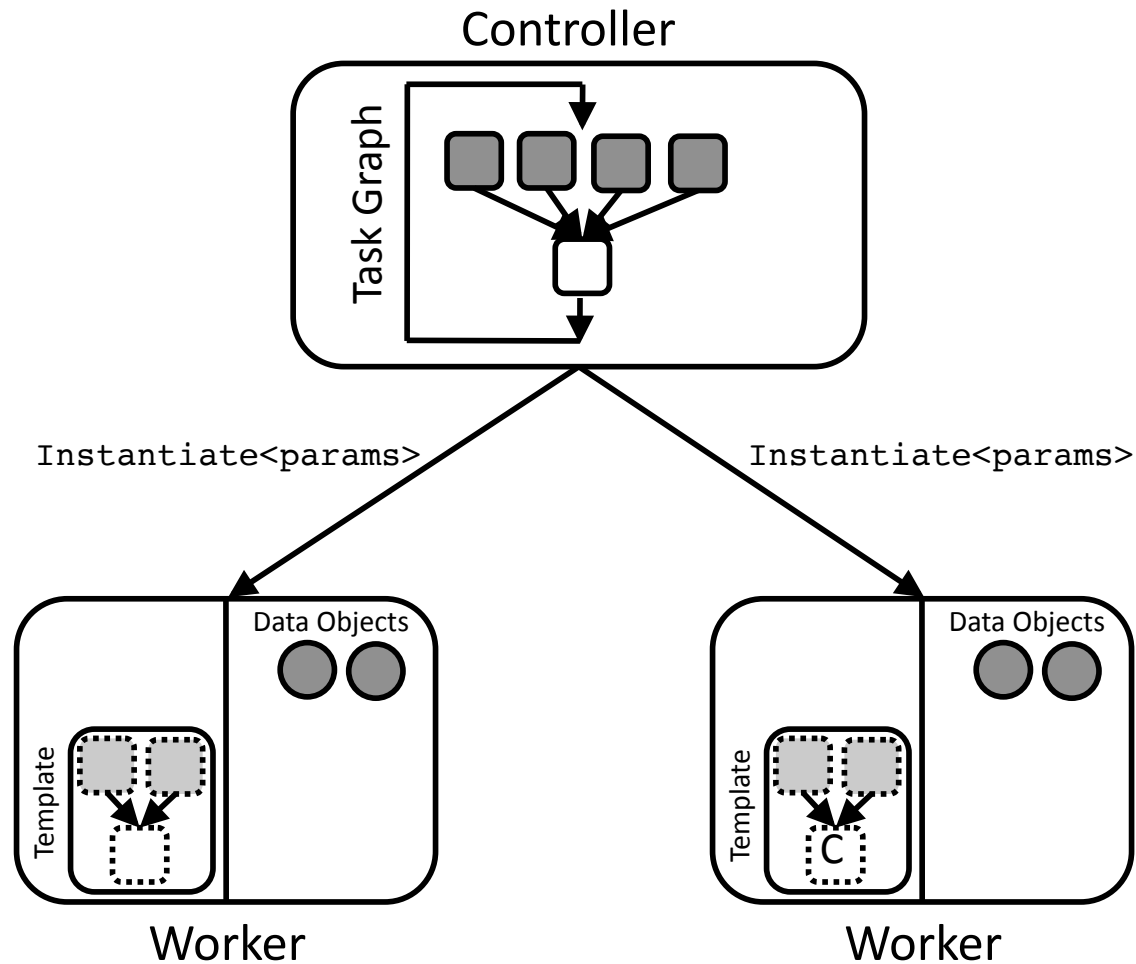
# Execution Templates

## Abstraction



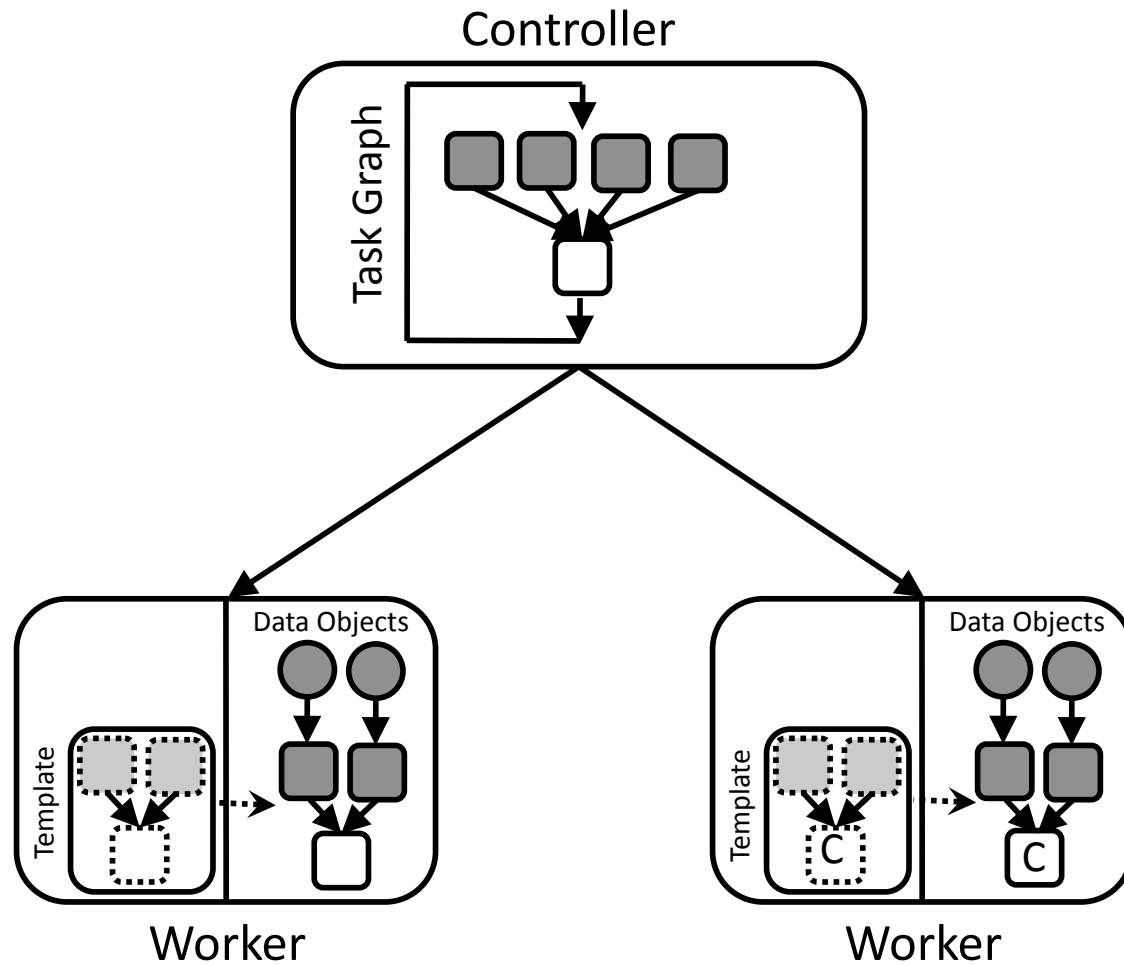
# Execution Templates

## Abstraction



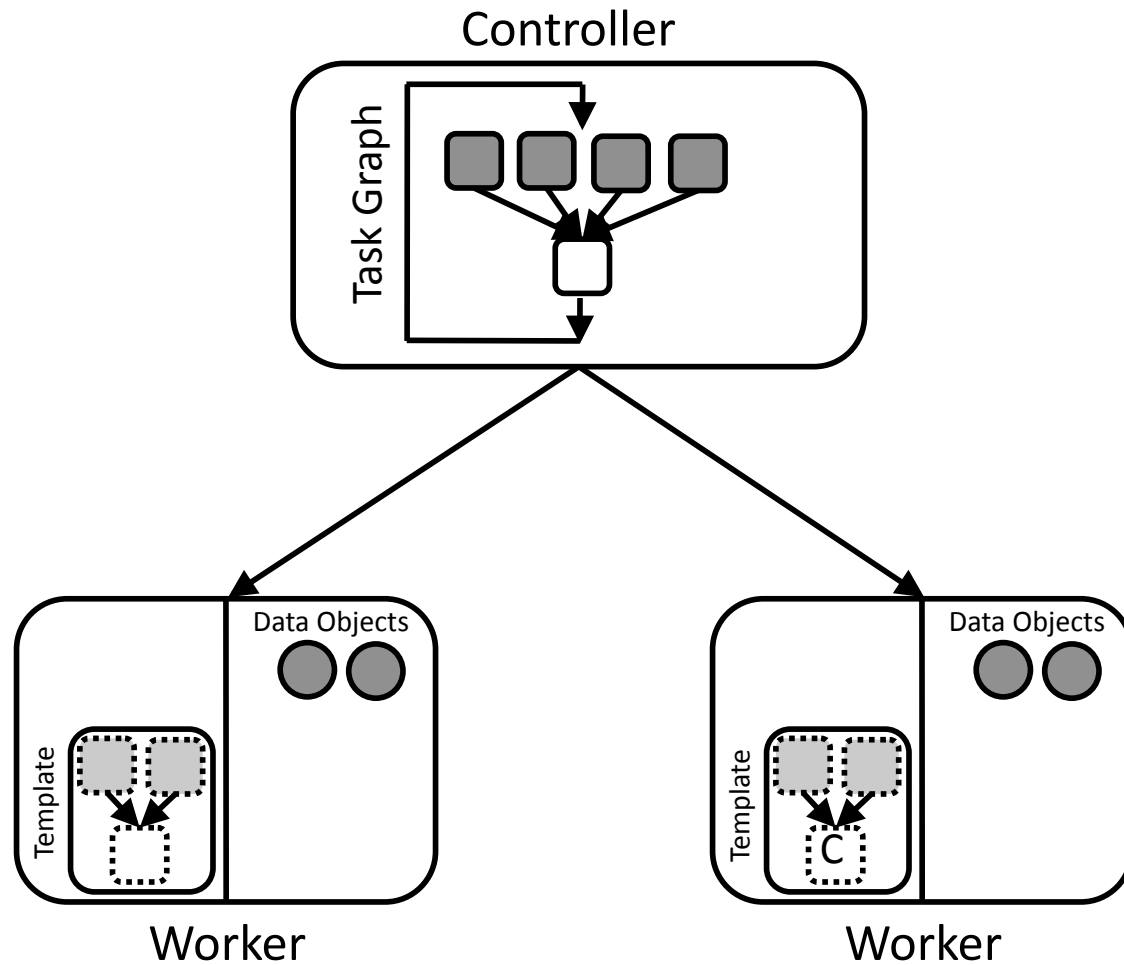
# Execution Templates

## Abstraction



# Execution Templates

## Abstraction



# Execution Templates

*The Devil is in the details.*

- Caching tasks implies static behavior:
  - Templates and **dynamic scheduling**?
    - Reactive scheduling changes for load balancing.
    - Scheduling changes at the task granularity.
  - Templates and **dynamic control flow**?
    - Need to support nested loops.
    - Need to support data dependent branches.

# Execution Templates

*The Devil is in the details.*

- Caching tasks implies static behavior:
  - Templates and **dynamic scheduling**?
    - Reactive scheduling changes for load balancing.
    - Scheduling changes at the task granularity.
  - Templates and **dynamic control flow**?
    - Need to support nested loops.
    - Need to support data dependent branches.

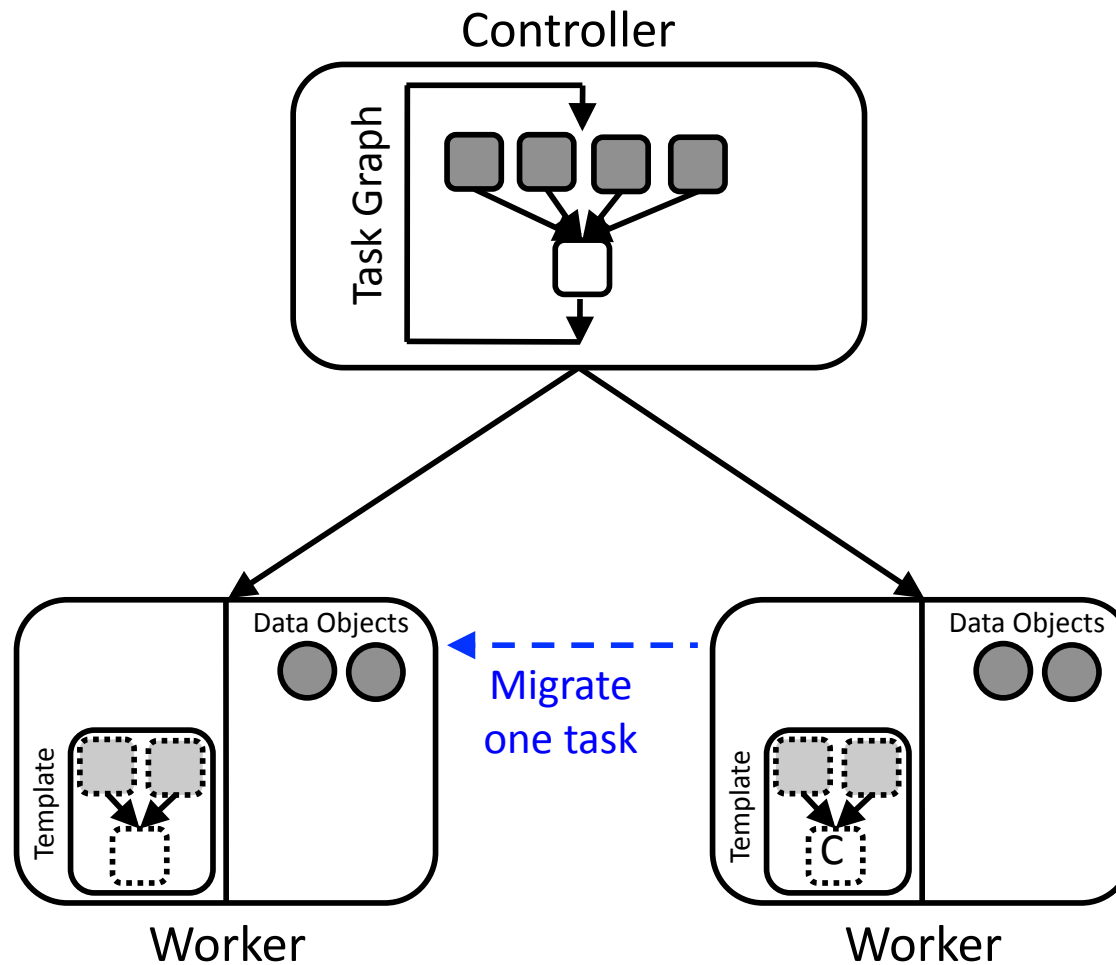
# Execution Templates

## Edits

- If scheduling changes, even slightly, the templates are obsolete.
  - For example migrating tasks among workers.
- Instead of paying the substantial cost of installing templates for every changes, templates allow **edit**, to change their structure.
- **Edits** enable adding or removing tasks from the template and modifying the template content, in-place.
- Controller has the general view of the task graph so it can update the dependencies properly, needed by the edits.

# Execution Templates

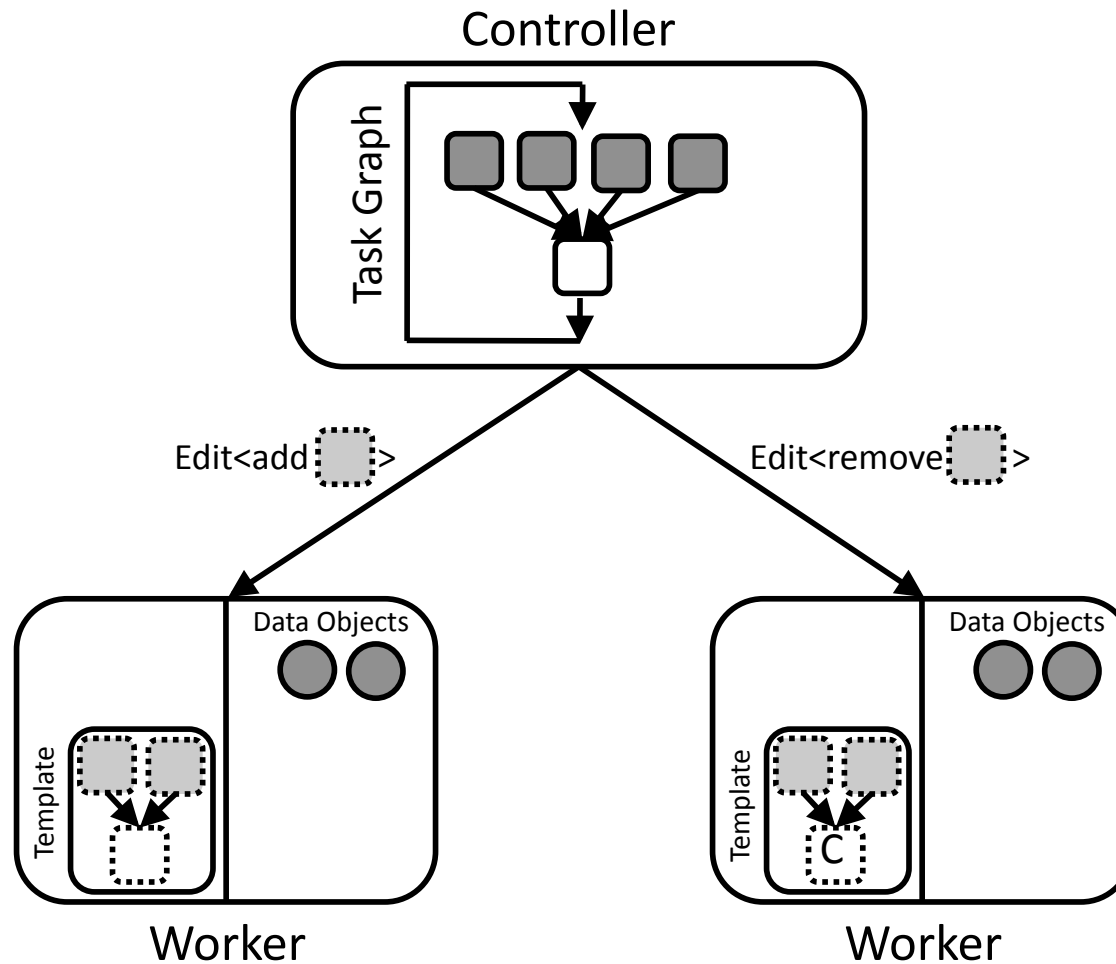
## Edits





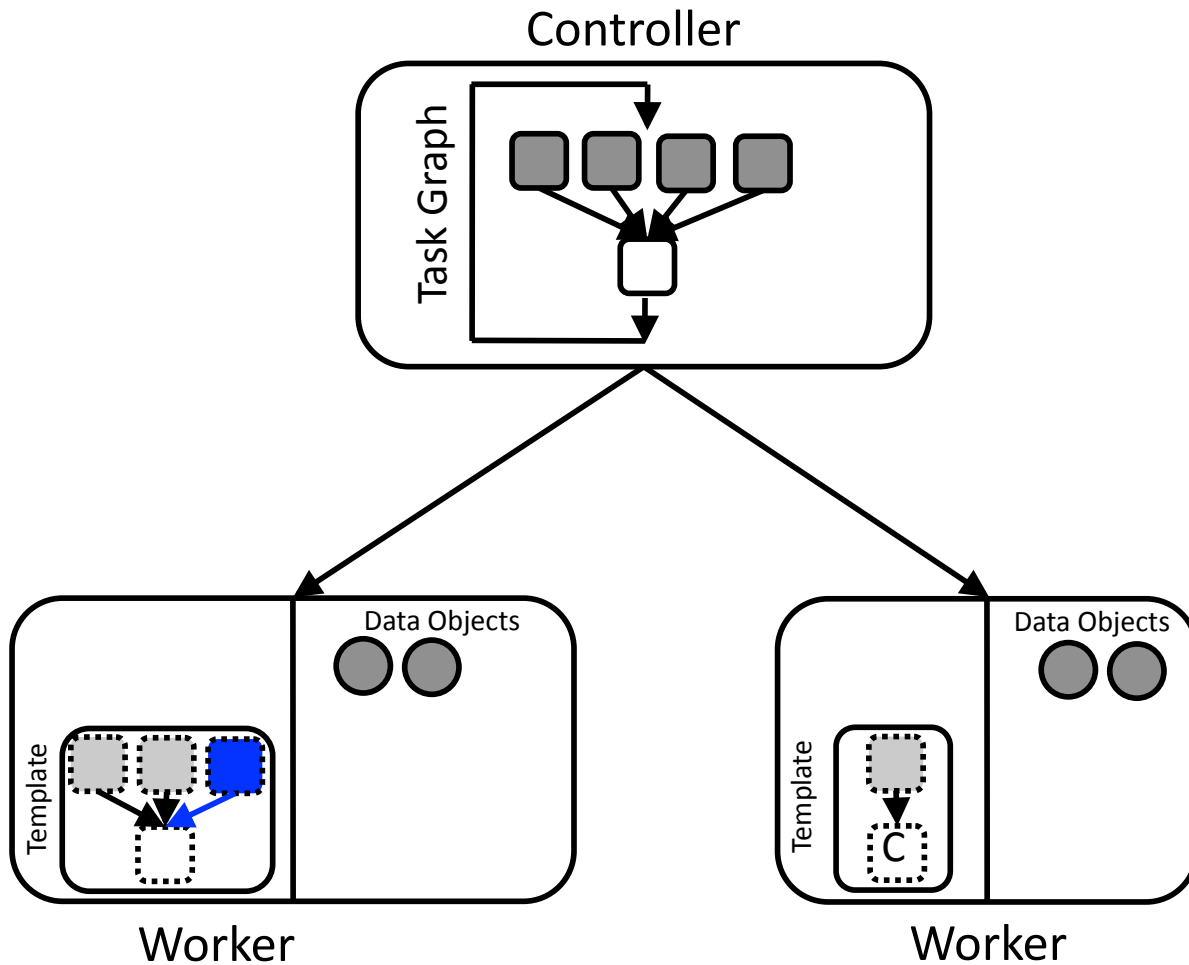
# Execution Templates

## Edits



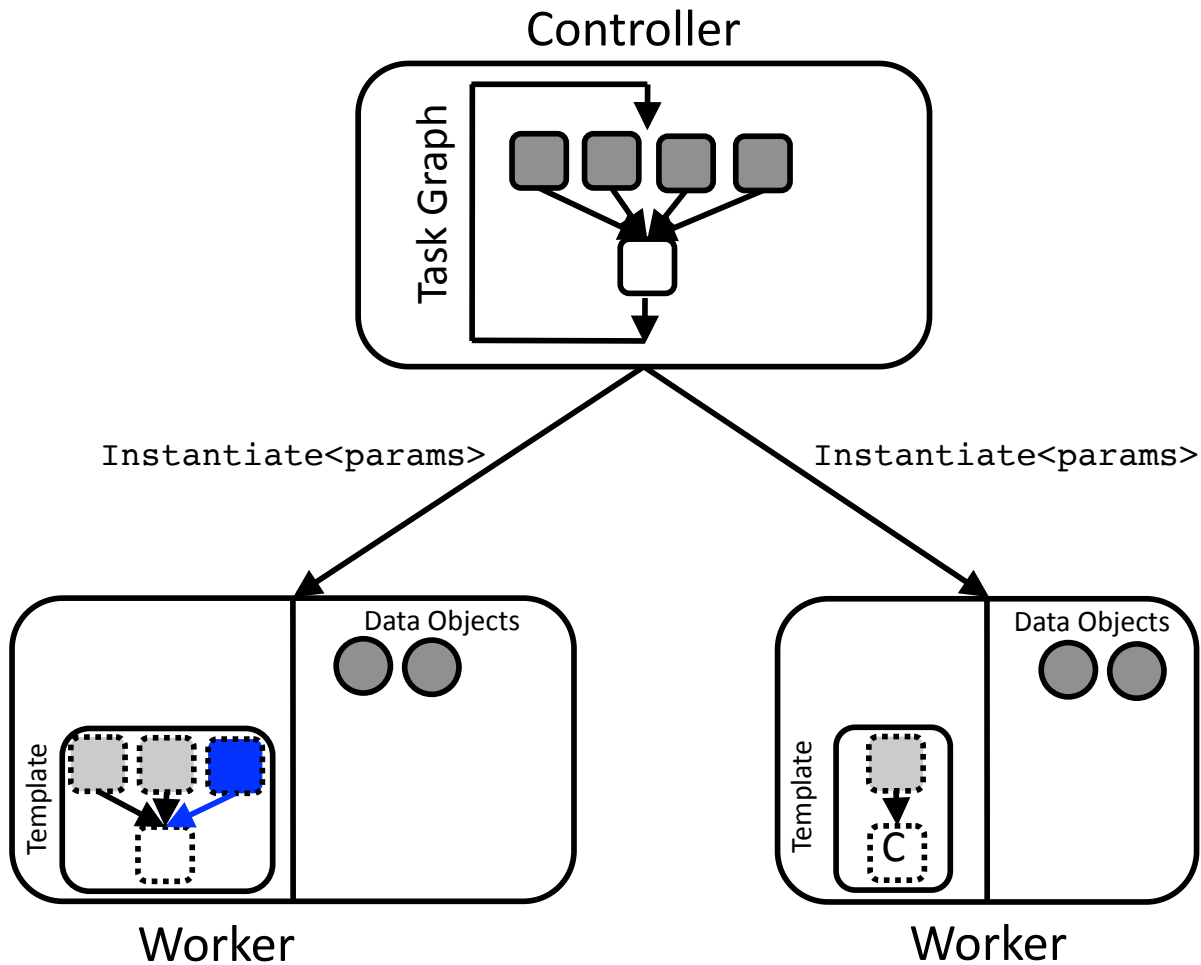
# Execution Templates

## Edits



# Execution Templates

## Edits



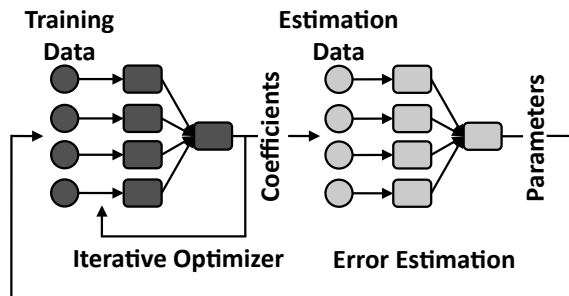
# Execution Templates

*The Devil is in the details.*

- Caching tasks implies static behavior:
  - Templates and **dynamic scheduling**?
    - Reactive scheduling changes for load balancing.
    - Scheduling changes at the task granularity.
  - Templates and **dynamic control flow**?
    - Need to support nested loops.
    - Need to support data dependent branches.

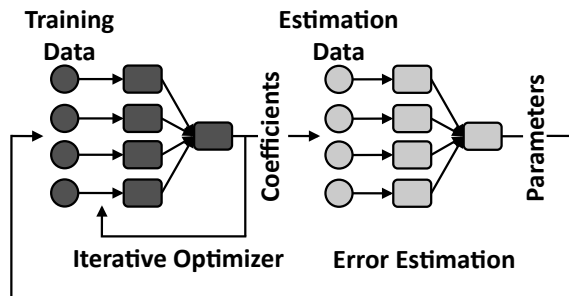
# Execution Templates

## Granularity



# Execution Templates

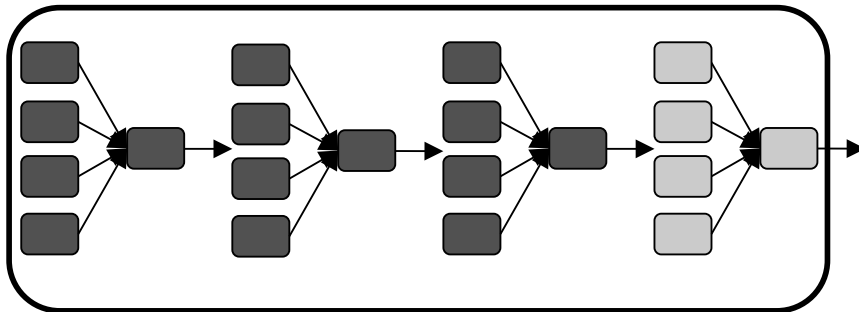
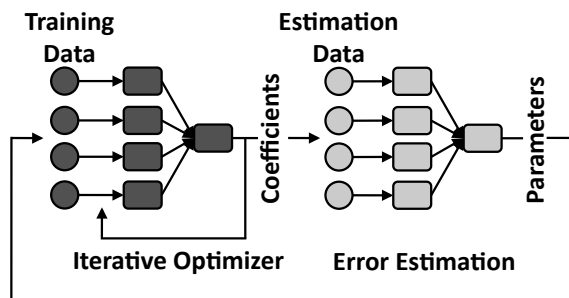
## Granularity



- The more tasks cached in the template the better.
  - The cost of template instantiation is amortized over greater number of tasks.
  - But **loop unrolling** only works for static control flow.

# Execution Templates

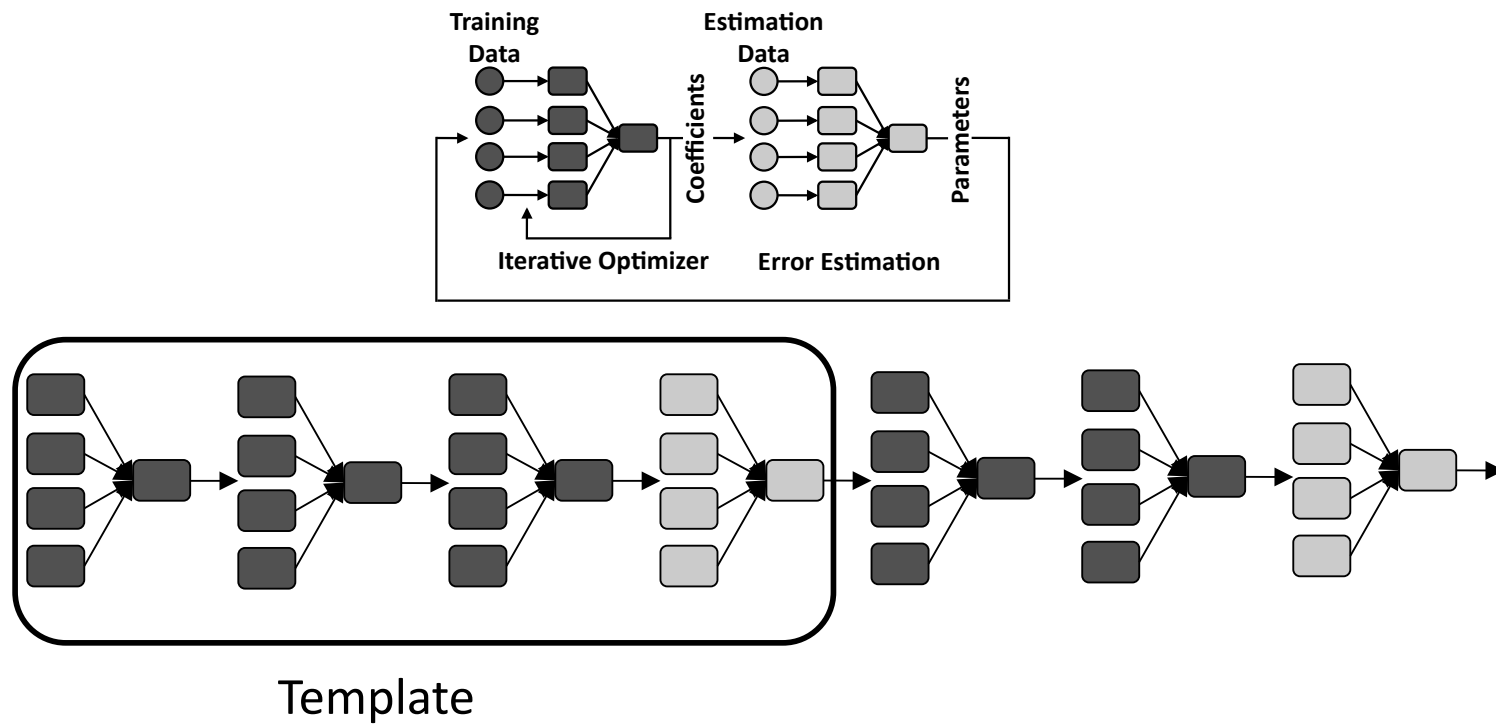
## Granularity



Template

# Execution Templates

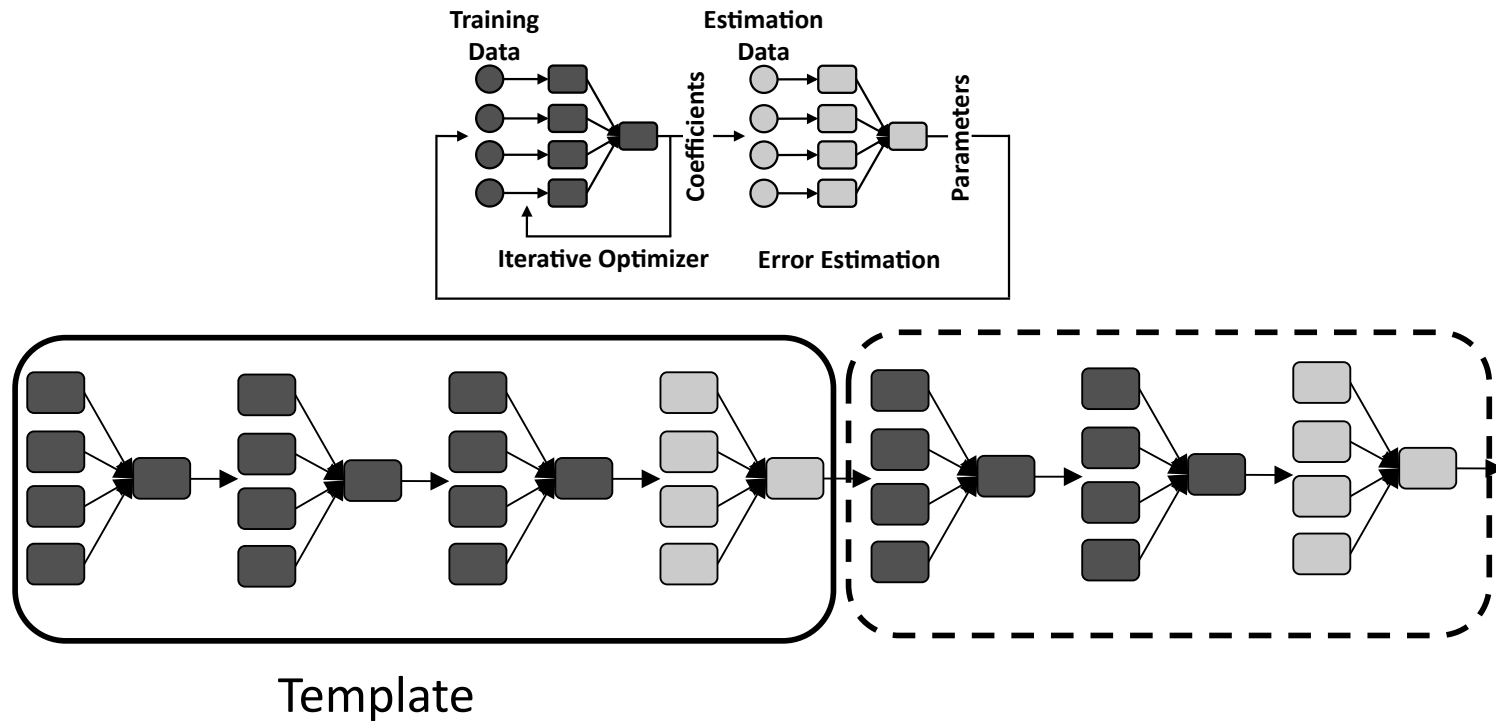
## Granularity





# Execution Templates

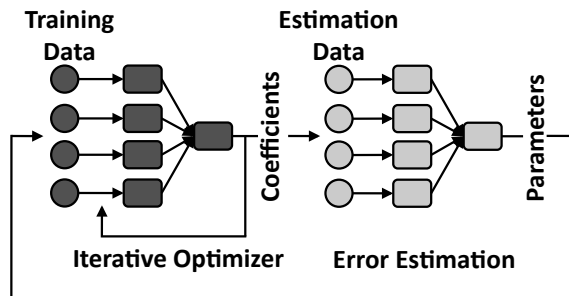
## Granularity



- Cannot reuse the template (only two iterations of the inner loop).

# Execution Templates

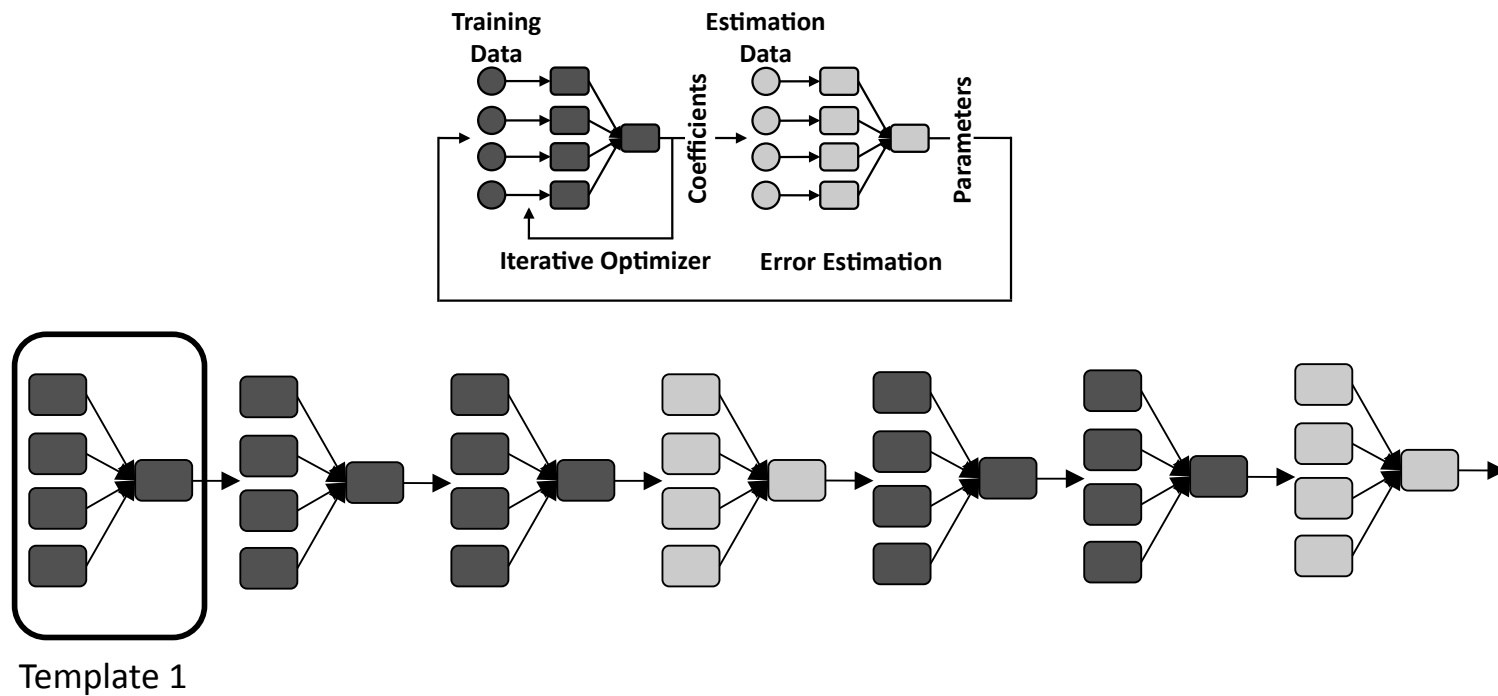
## Granularity



- Templates cannot go beyond a branch in the driver program.
- Execution templates operates at the granularity of **basic blocks**:
  - A code block with single entry and no branches except at the end.
  - It is the biggest block without sacrificing **dynamic control flow**.

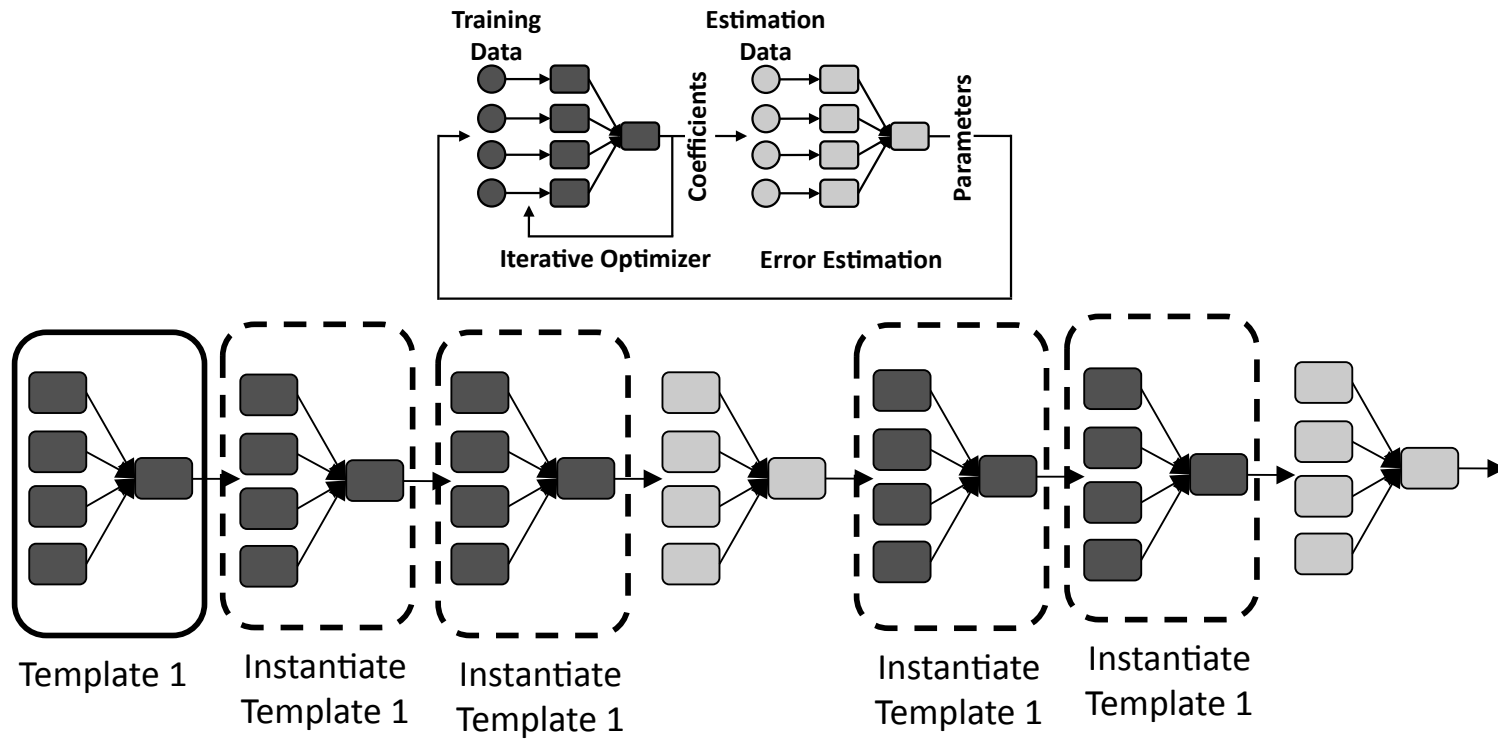
# Execution Templates

## Granularity



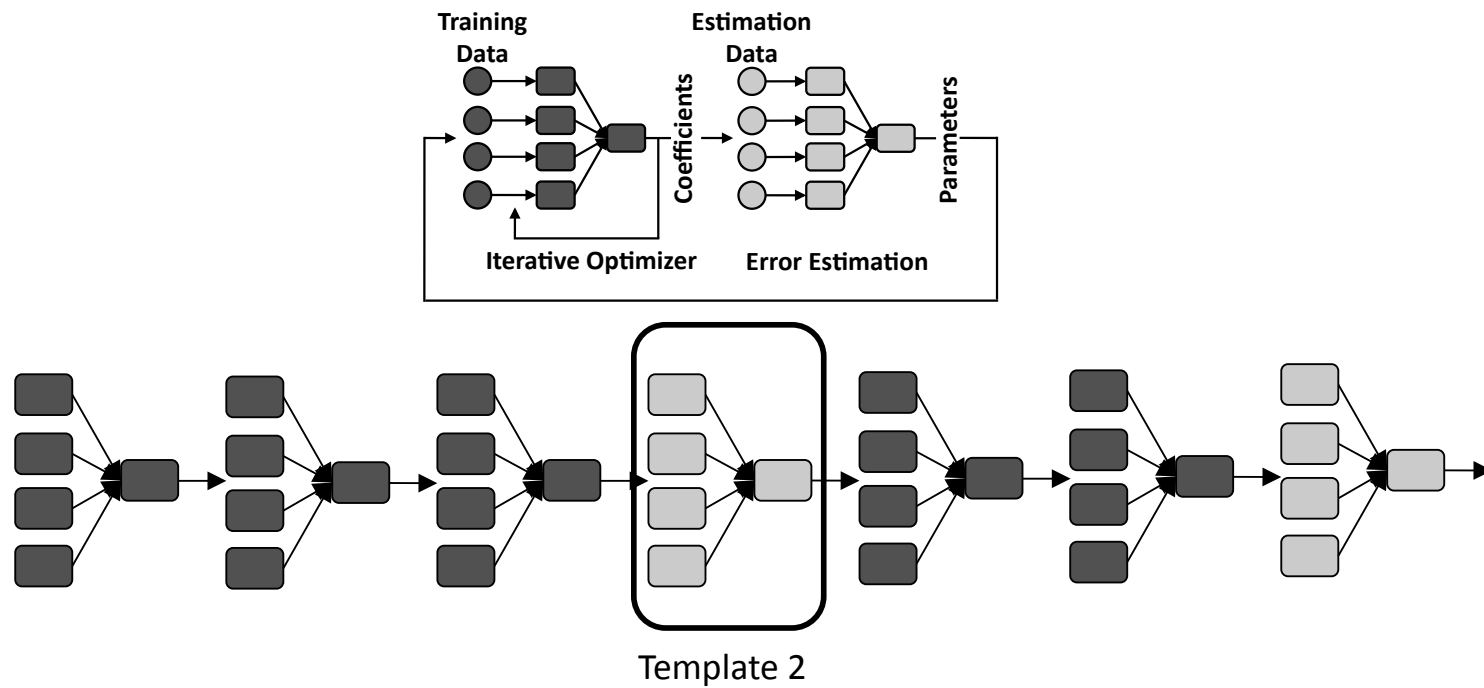
# Execution Templates

## Granularity



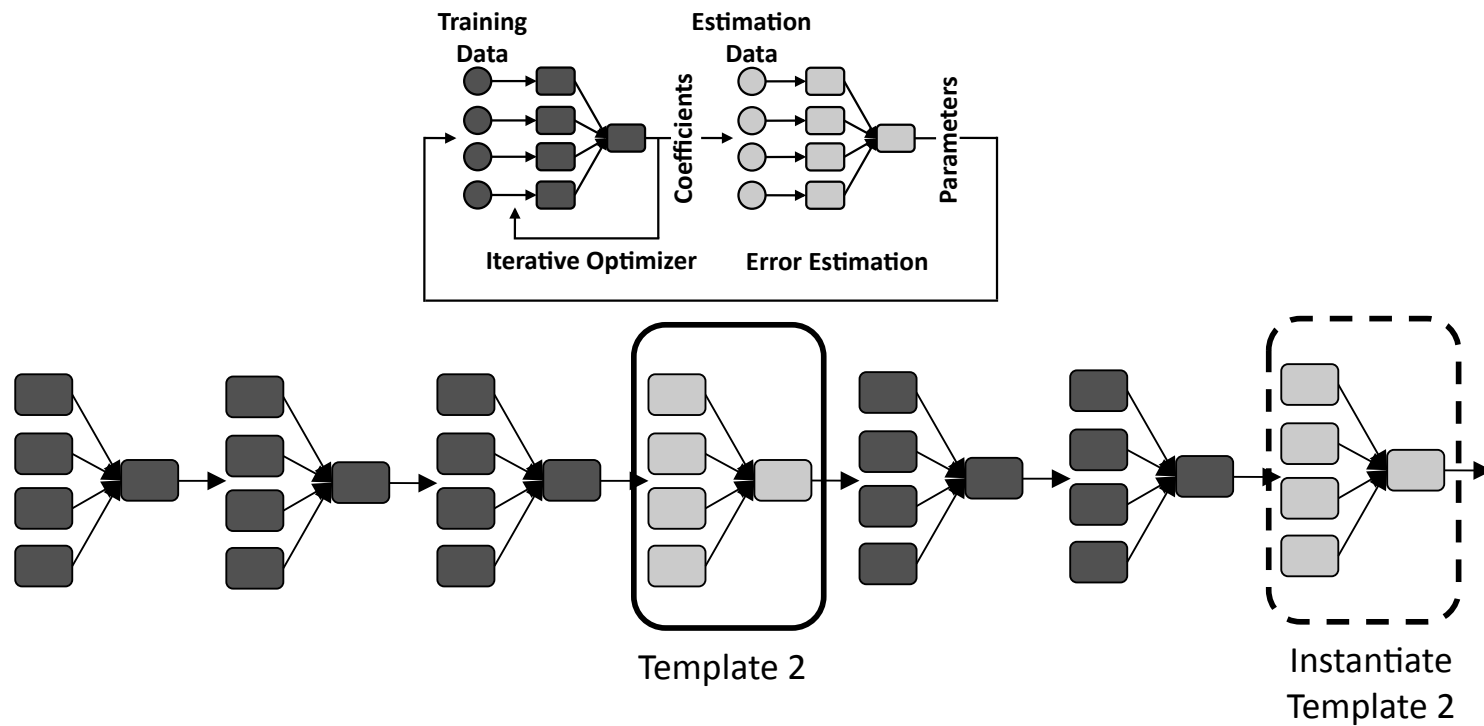
# Execution Templates

## Granularity



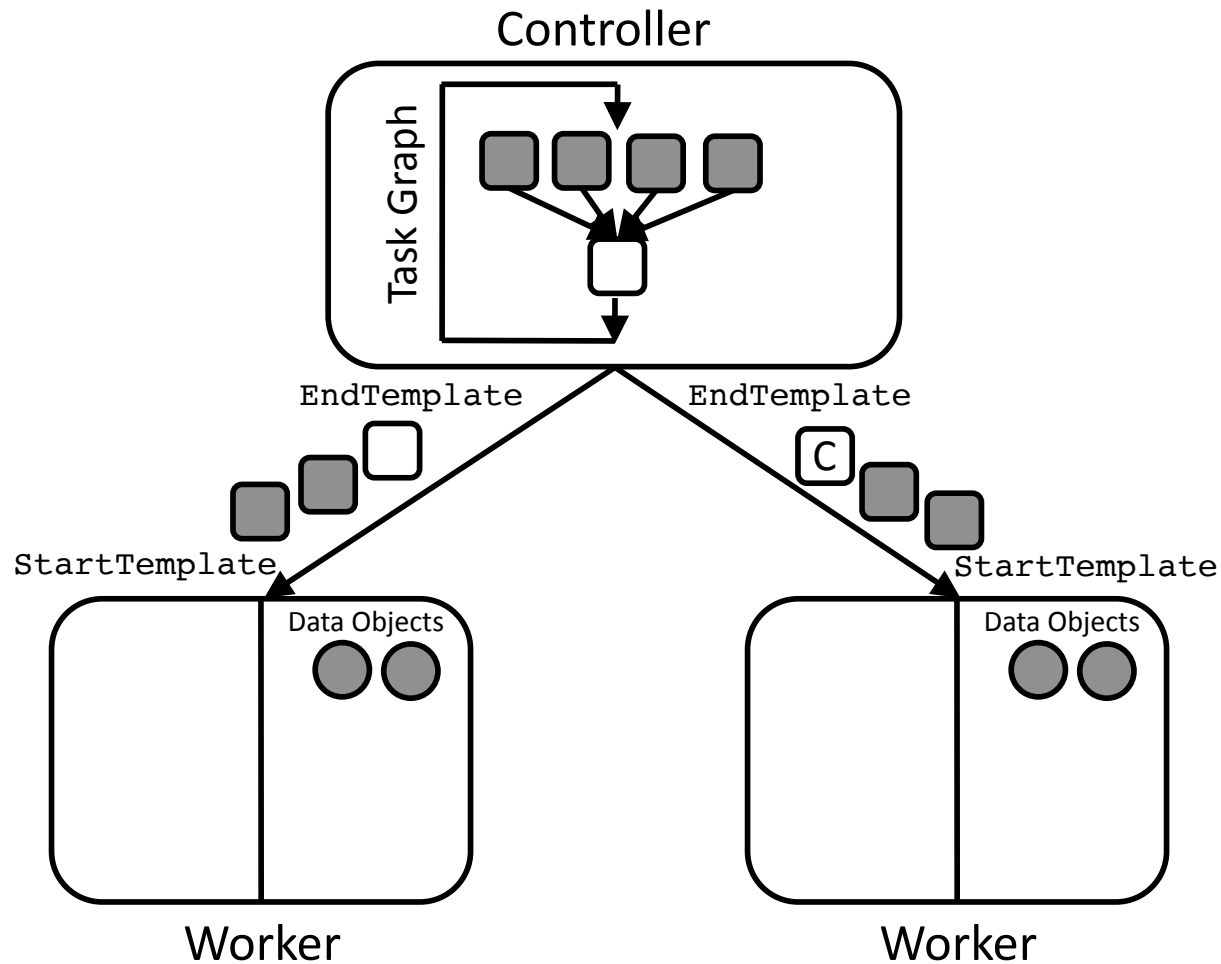
# Execution Templates

## Granularity



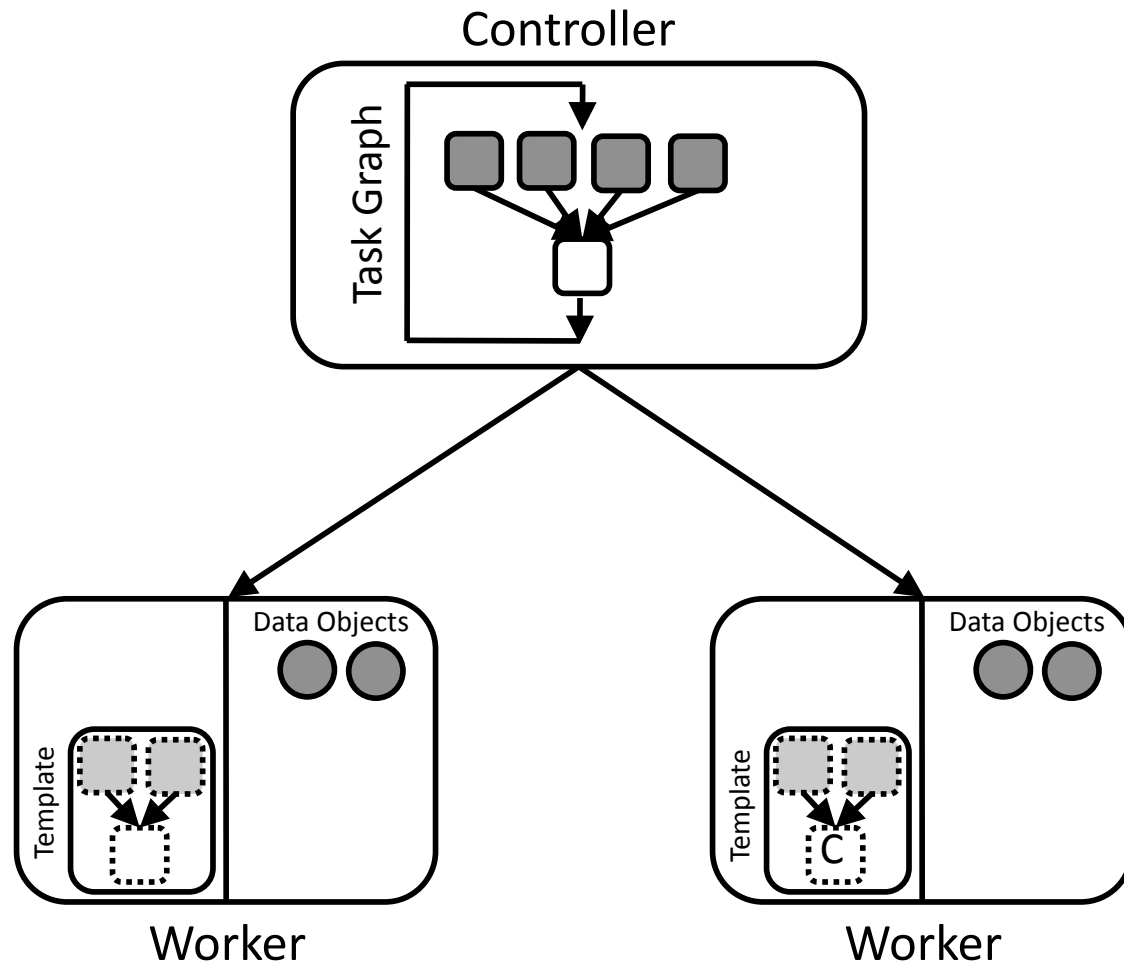
# Execution Templates

## Granularity



# Execution Templates

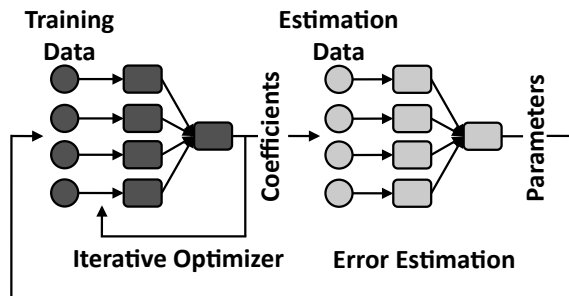
## Granularity





# Execution Templates

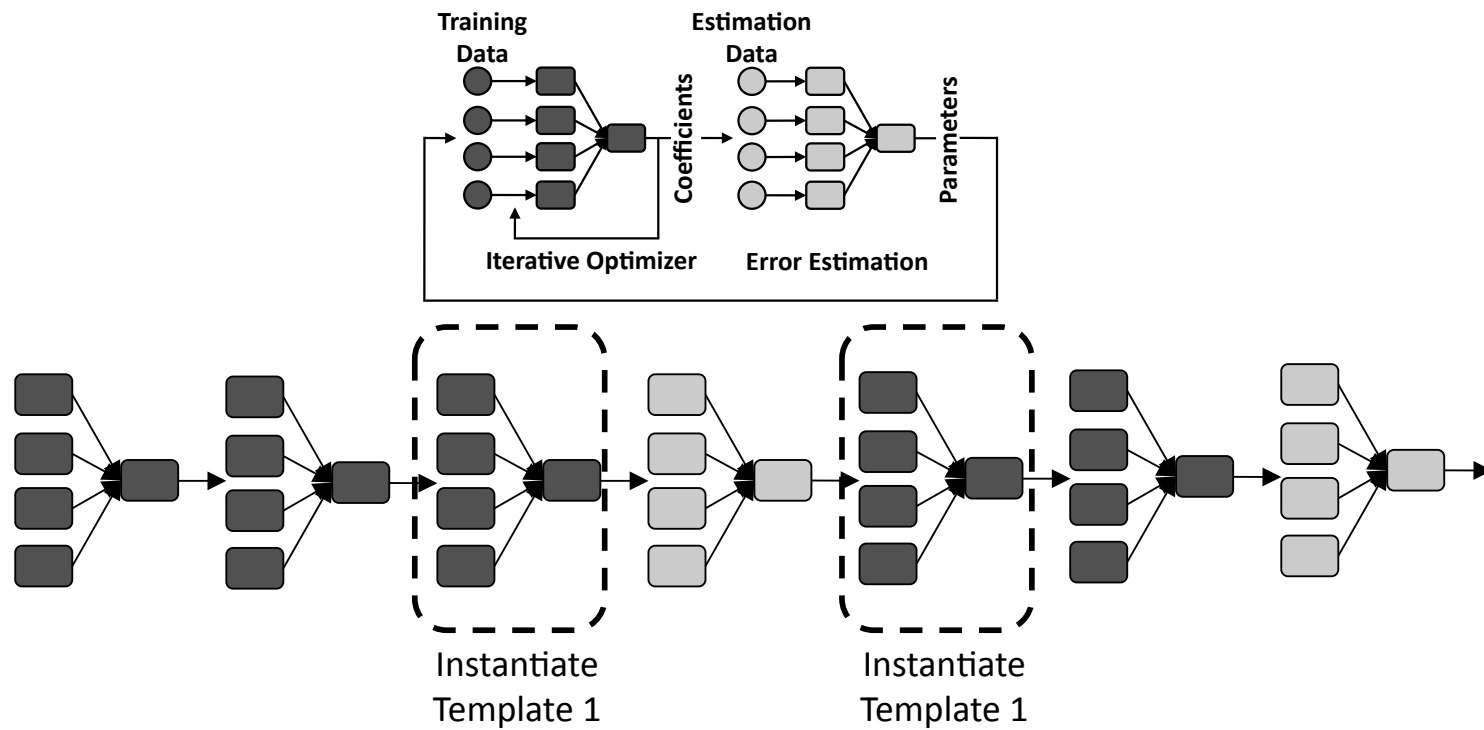
## Patching



- With dynamic control flow a basic block can have different entries.
- The execution state is not similar in all circumstances.

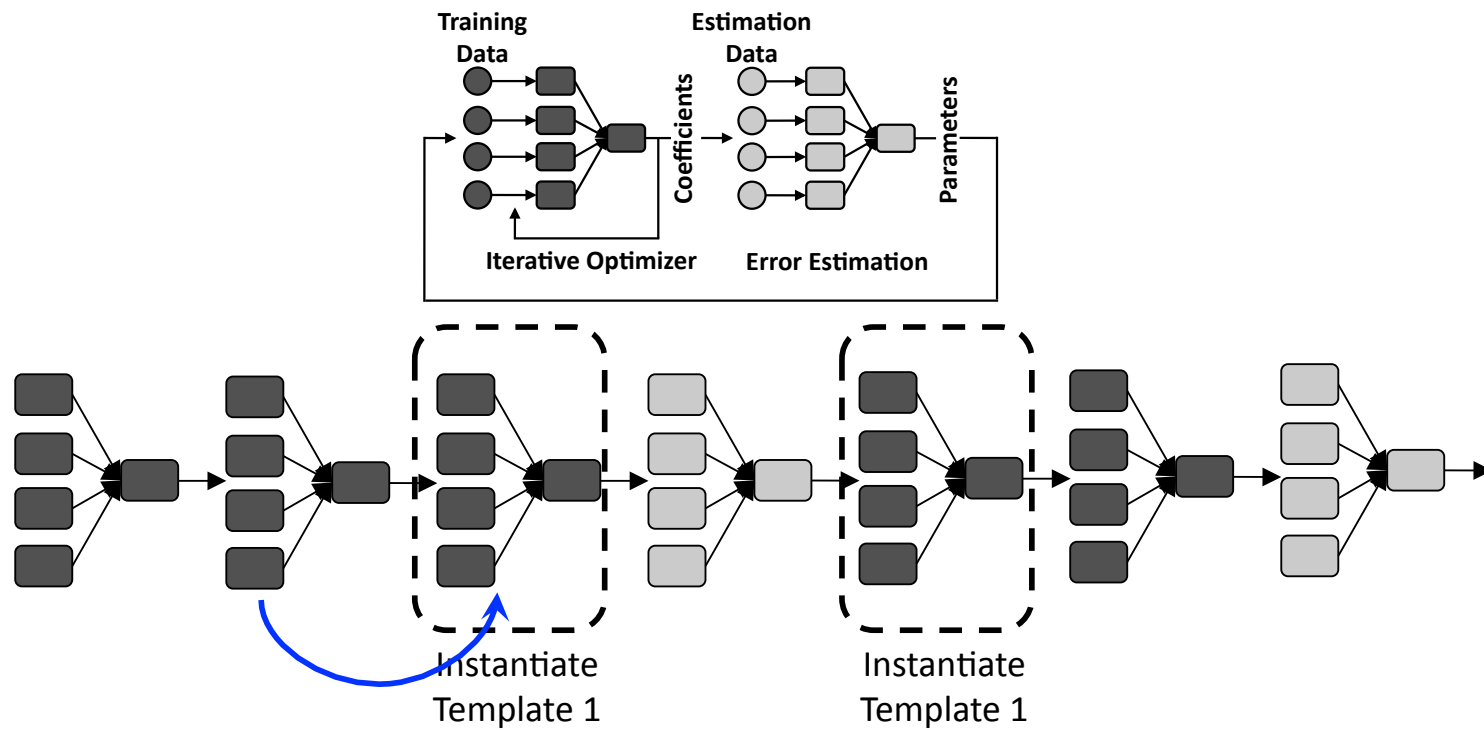
# Execution Templates

## Patching



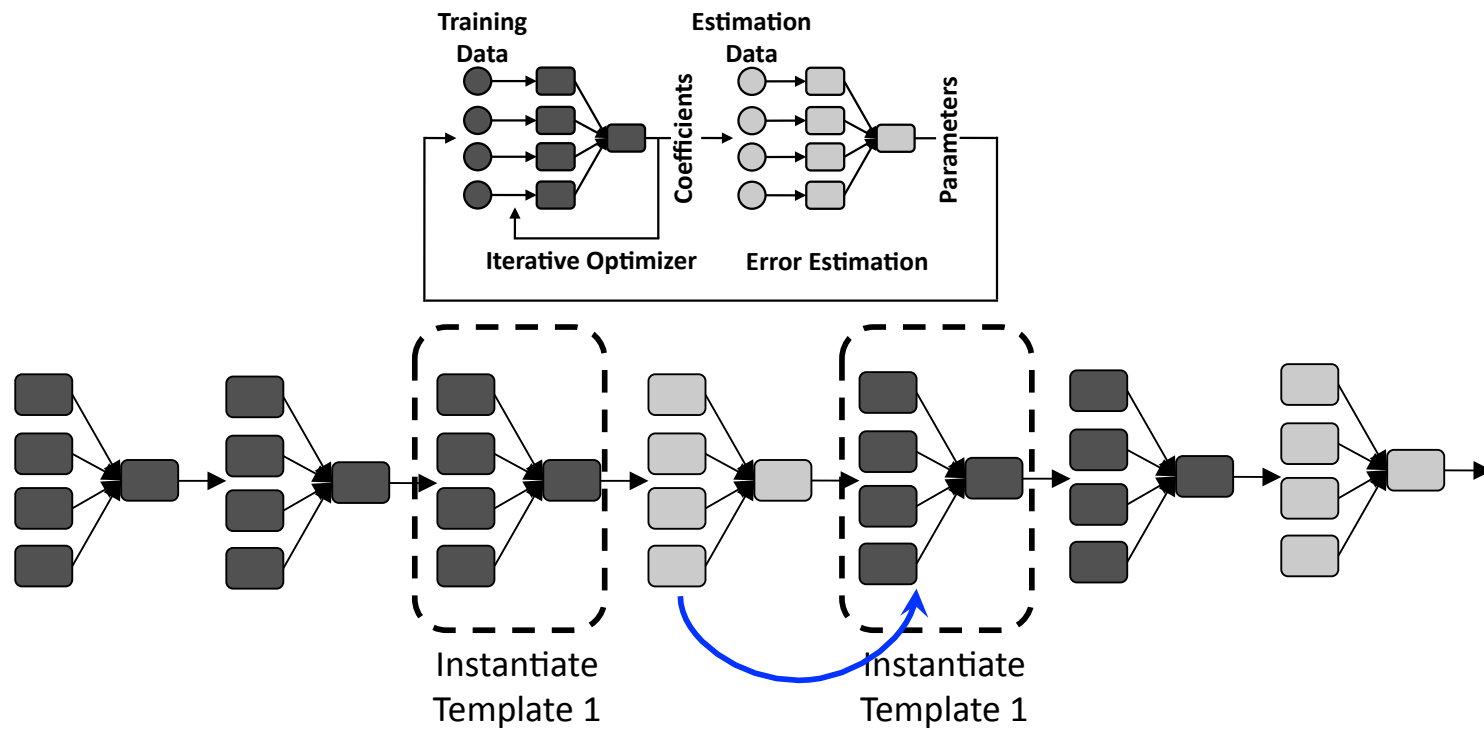
# Execution Templates

## Patching



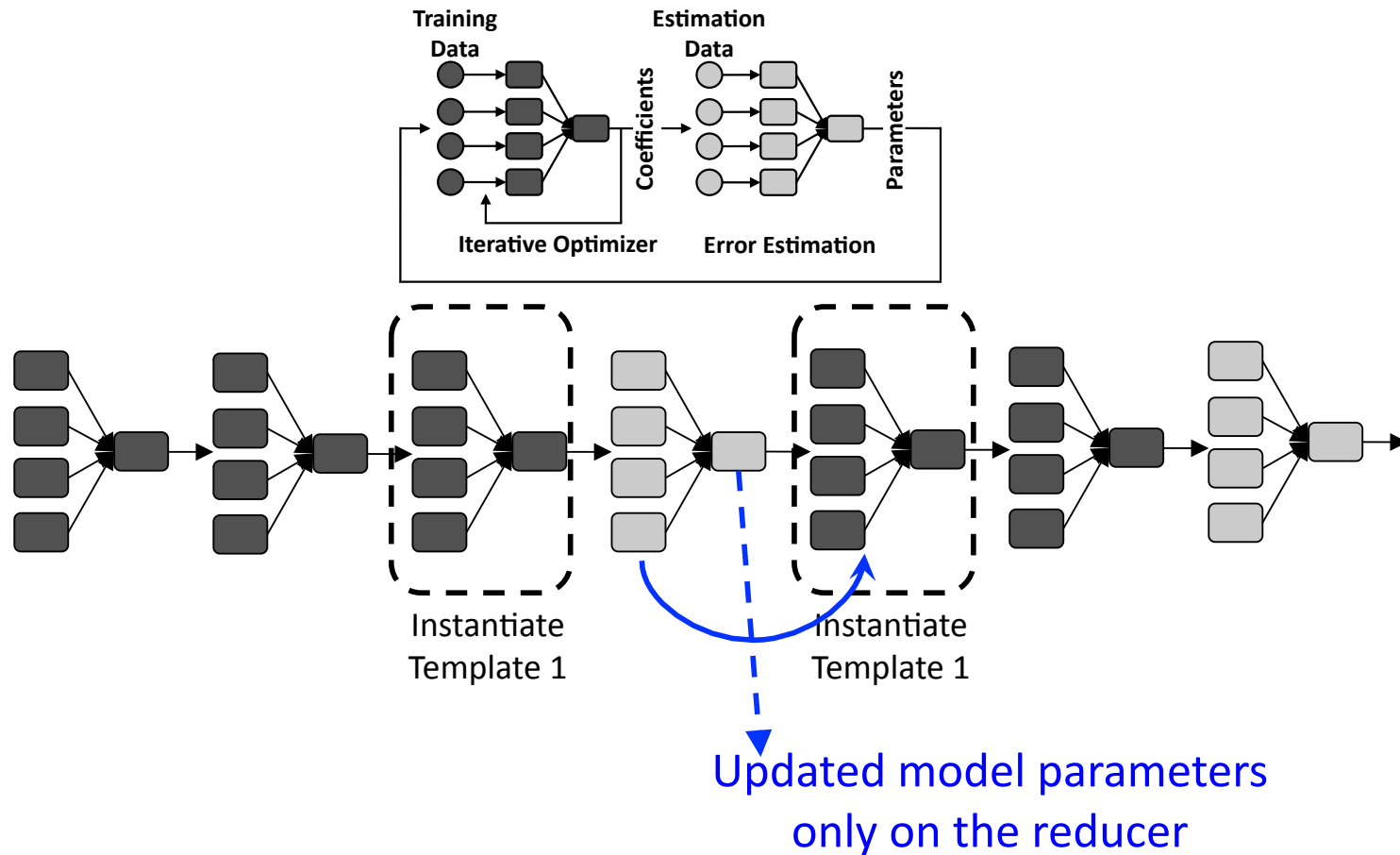
# Execution Templates

## Patching



# Execution Templates

## Patching



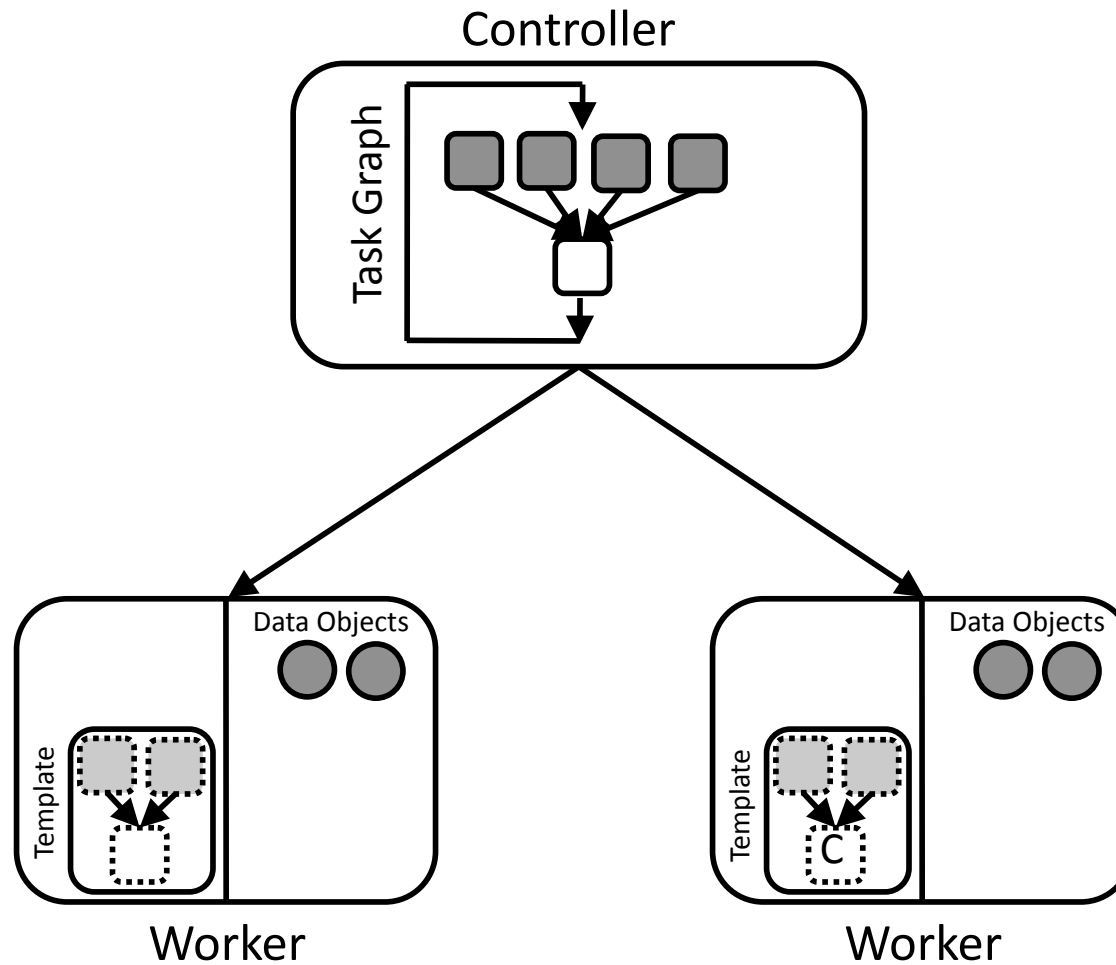
# Execution Templates

## Patching

- Each template has a set of **preconditions** that need to be satisfied before it can be instantiated.
  - For example the set of data objects in memory, accessed by the tasks cached in the template.

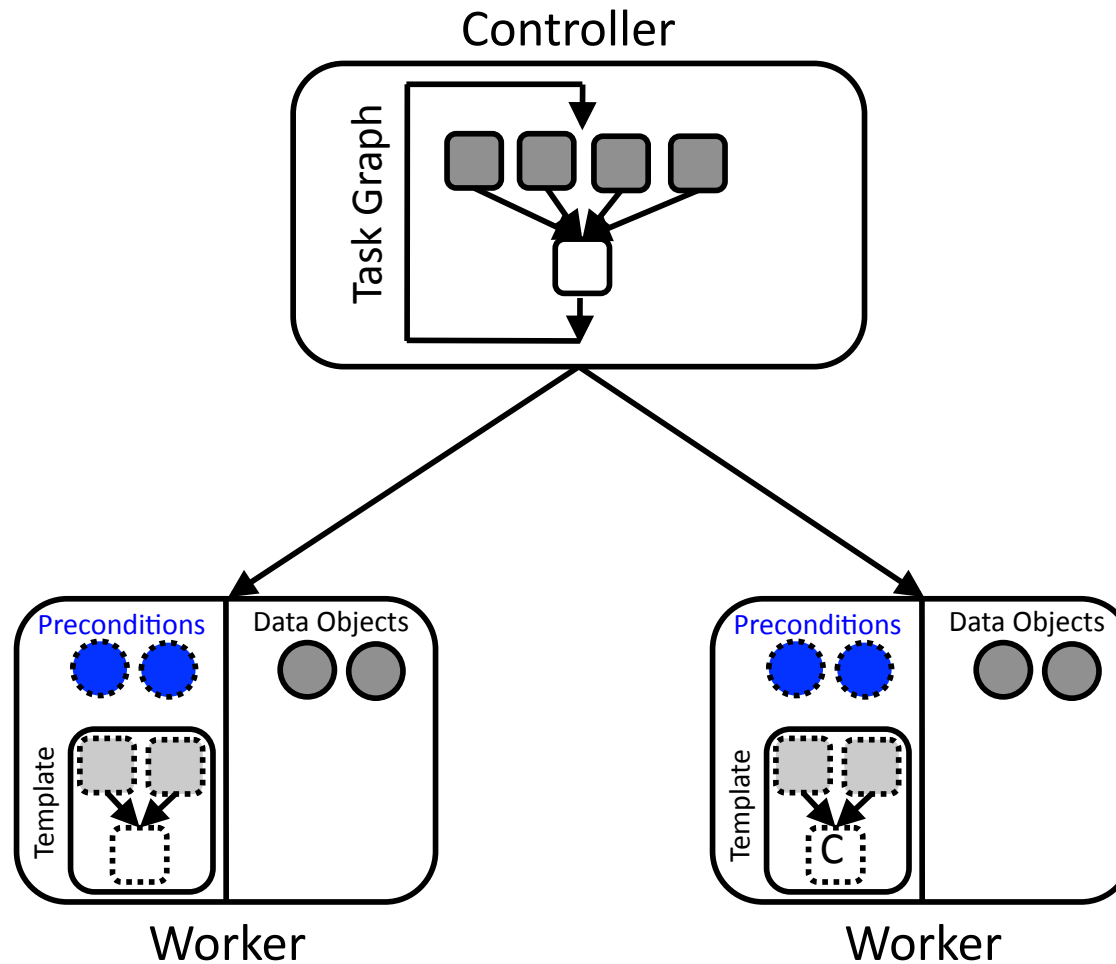
# Execution Templates

## Patching



# Execution Templates

## Patching





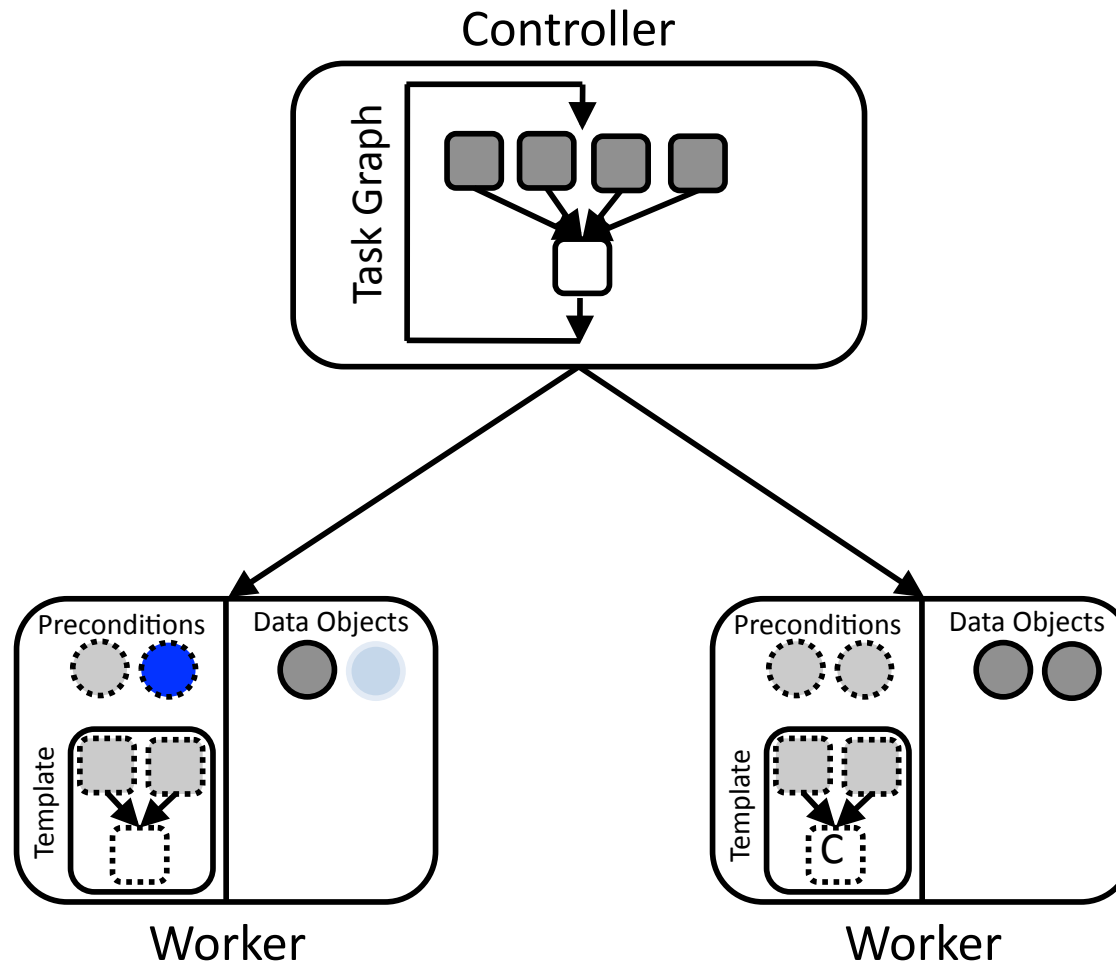
# Execution Templates

## Patching

- Each template has a set of **preconditions** that need to be satisfied before it can be instantiated.
  - For example the set of data objects in memory, accessed by the tasks cached in the template.
- Worker state might not match the preconditions of the template in all circumstances.
- Controller **patches** the worker state before template instantiation, to satisfy the preconditions.

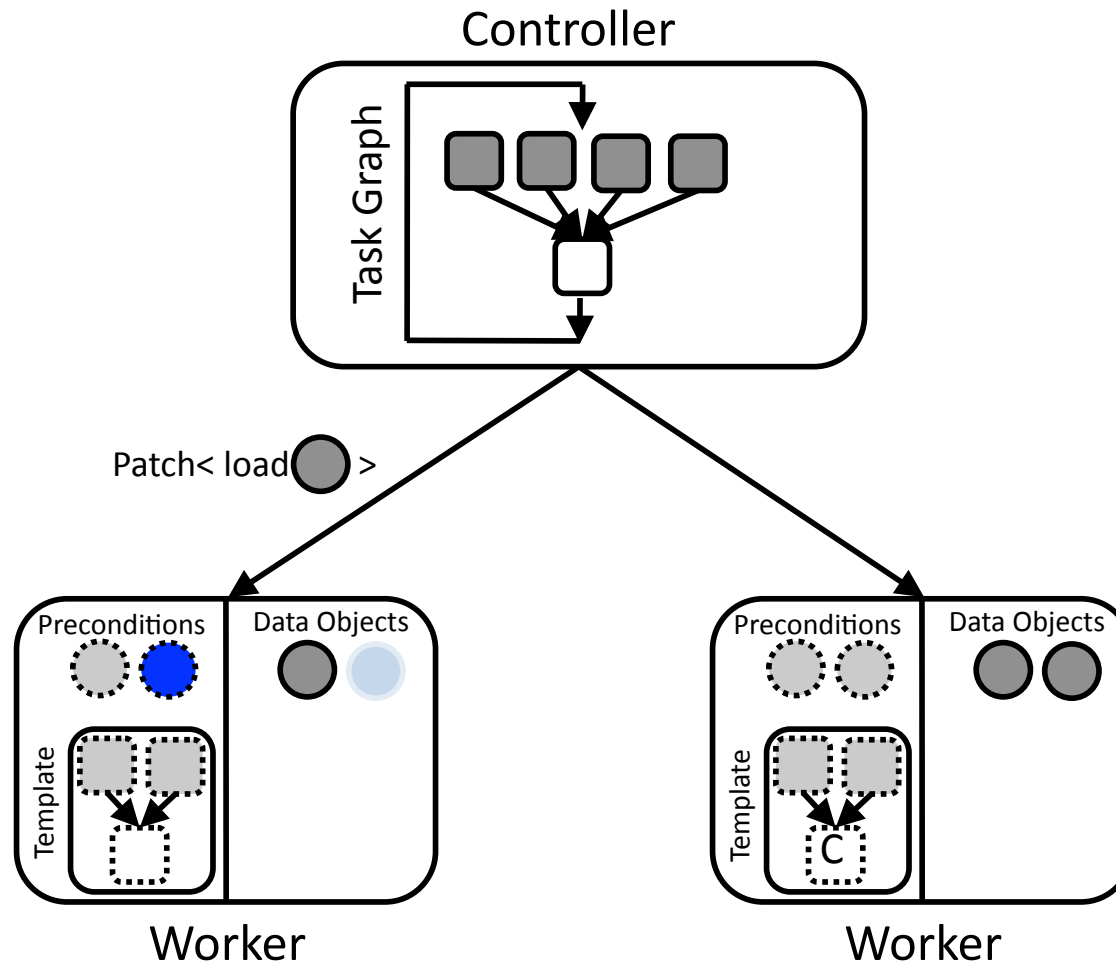
# Execution Templates

## Patching



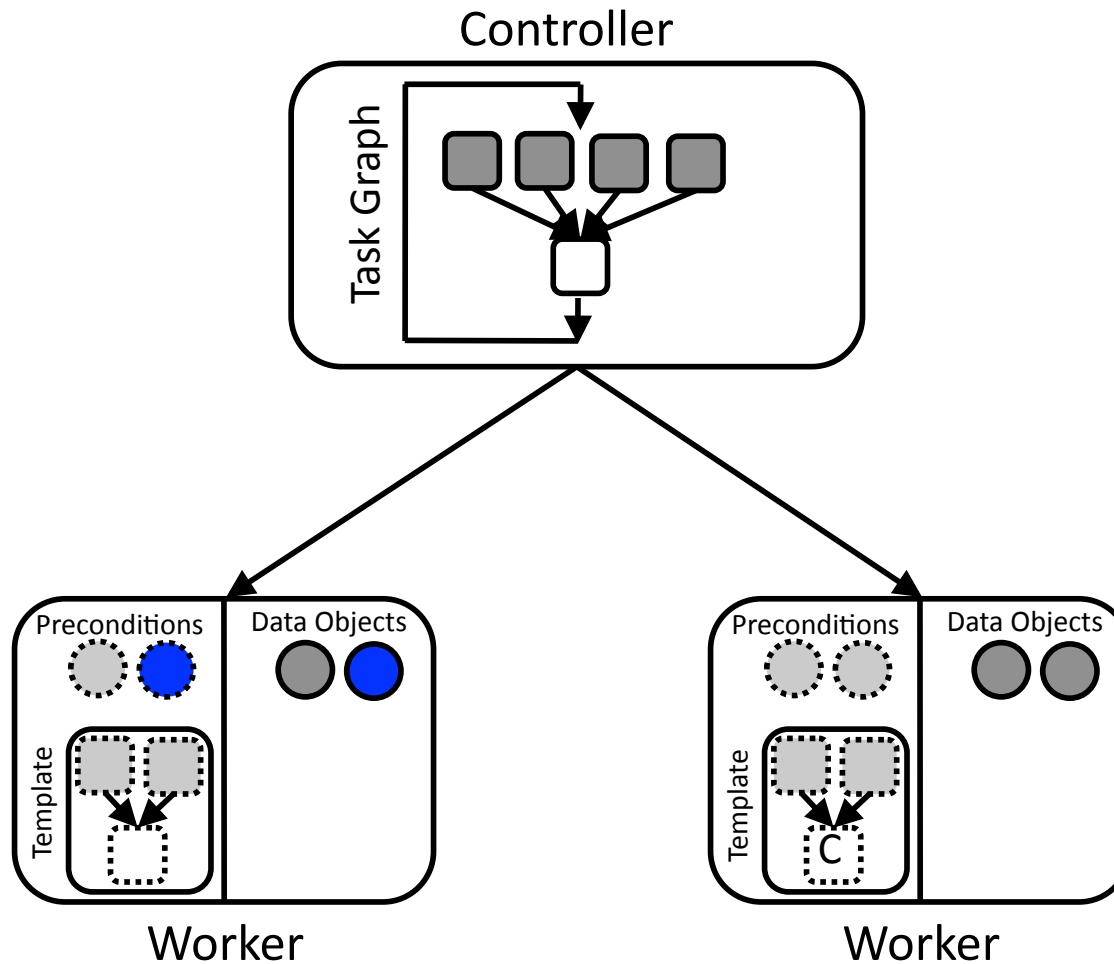
# Execution Templates

## Patching



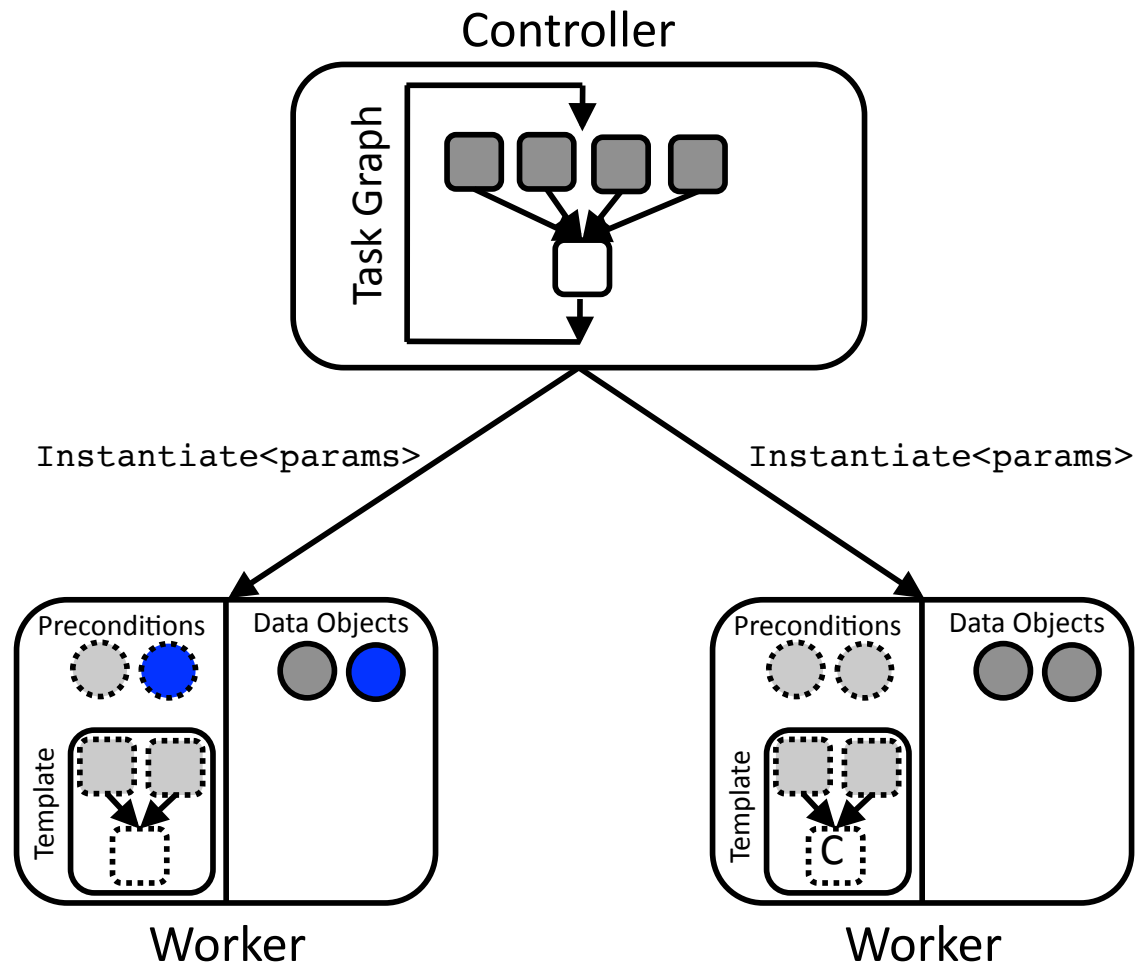
# Execution Templates

## Patching



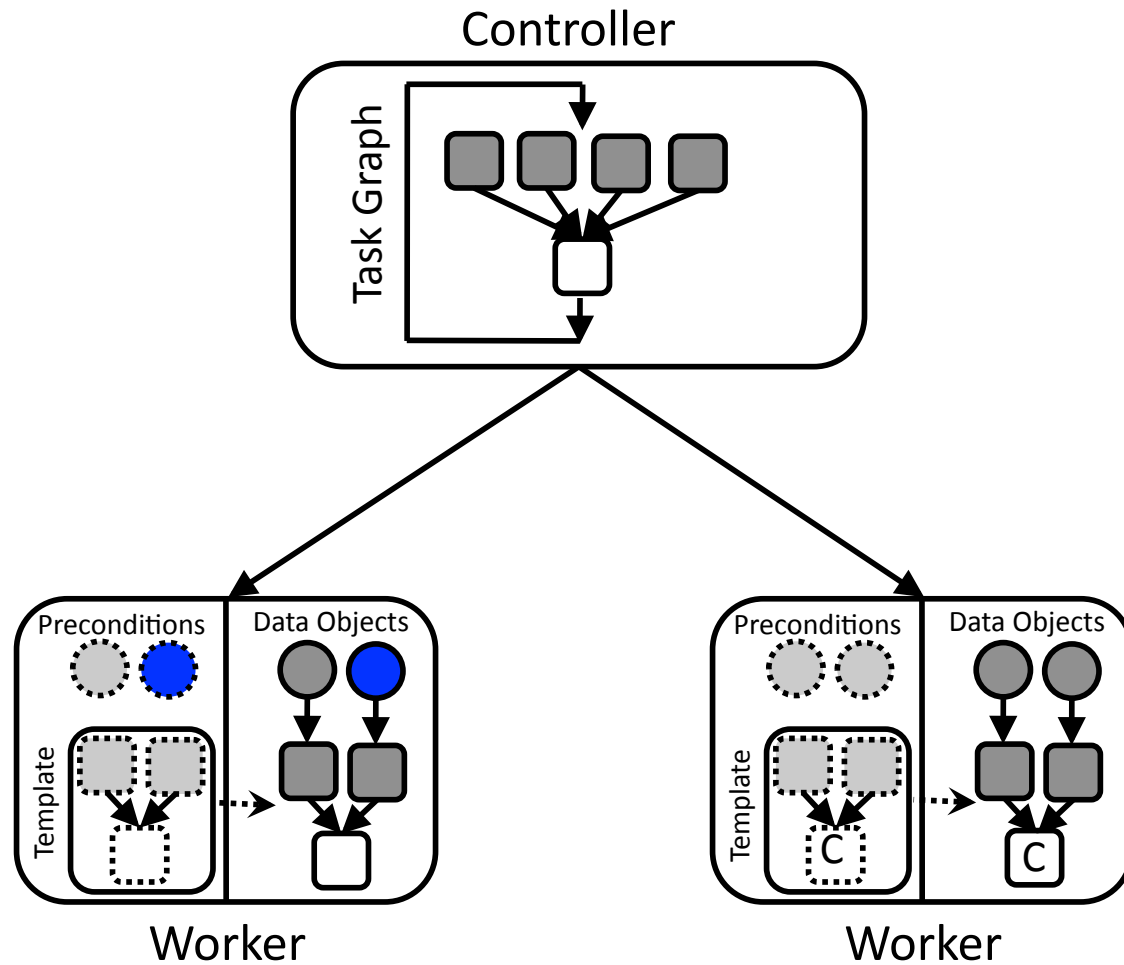
# Execution Templates

## Patching



# Execution Templates

## Patching



# Execution Templates

## Mechanisms Summary

- **Instantiation:** spawn a block of tasks without processing each task individually from scratch. It helps increase the **task throughput**.
- **Edits:** modifies the content of each template at the granularity of tasks. It enables fine-grained, **dynamic scheduling**.
- **Patches:** In case the state of the worker does not match the preconditions of the template. It enables **dynamic control flow**.

# This talk

- Control Plane: the Emerging Bottleneck
- Design Scope of the Control Plane
- Execution Templates
- Nimbus: a Framework with Templates
- Evaluation

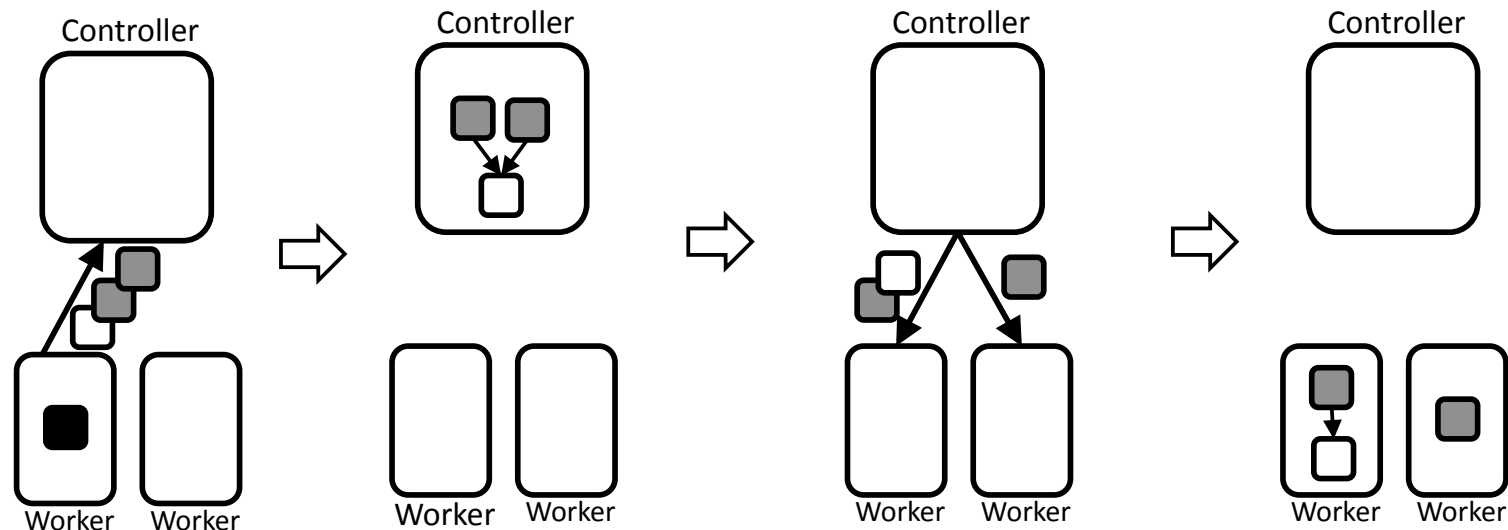


# Nimbus

- Nimbus is designed for low latency, fast computations in the cloud.
  - Implemented in C++ (the core library is ~35,000 semicolons).
  - Mutable data model to allow in-place operations.
- Nimbus embeds execution templates for its control plane.
  - The centralized controller allows dynamic scheduling and resource allocation.
  - Execution templates help deliver high task throughput at scale.
- Nimbus supports traditional data analytics as well as Eulerian and hybrid graphical simulations; for the first time in a cloud framework.
  - Supervised/unsupervised learning algorithms, graph library.
  - PhysBAM library (water, smoke, etc.)

# Nimbus

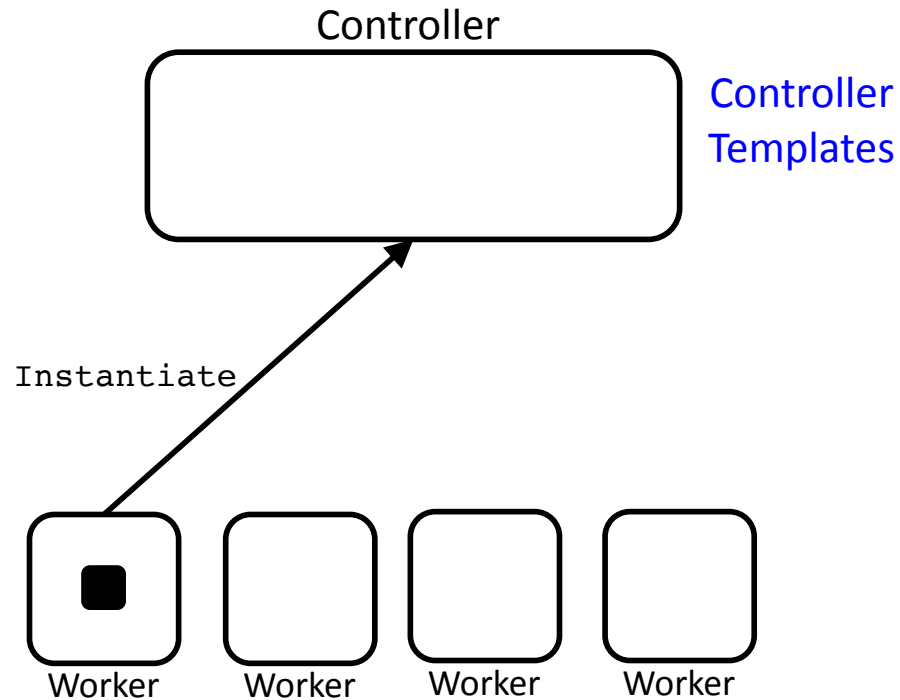
## Control Flow



- Tasks spawn other tasks for execution (similar to Legion).
- Driver program is a lineage of tasks executing on the workers.
- More flexible DAG for the task graph.
  - Not just narrow and wide dependencies.
  - Needed for graphical simulations.

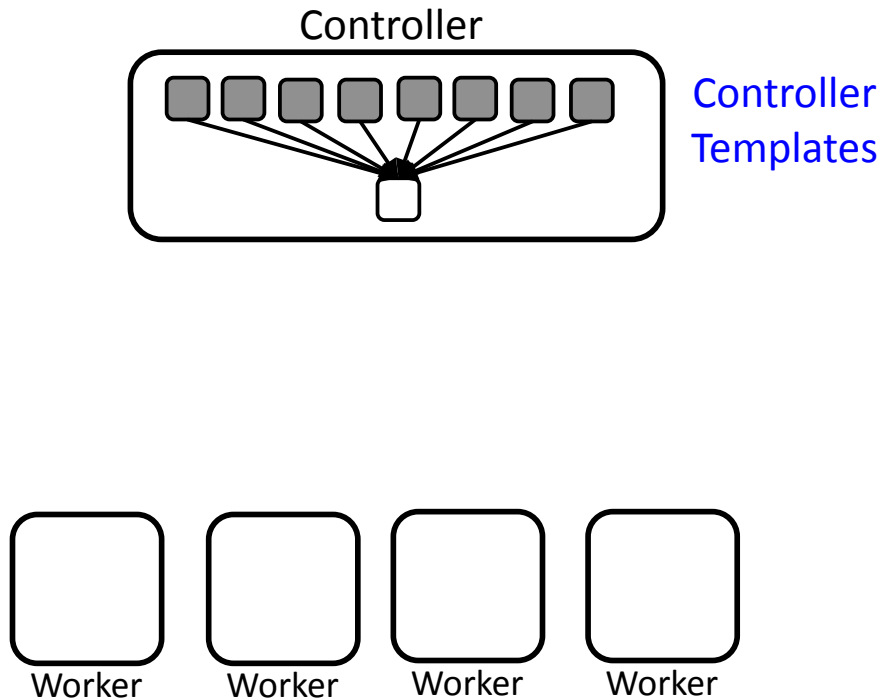
# Nimbus

## Controller and Worker Templates



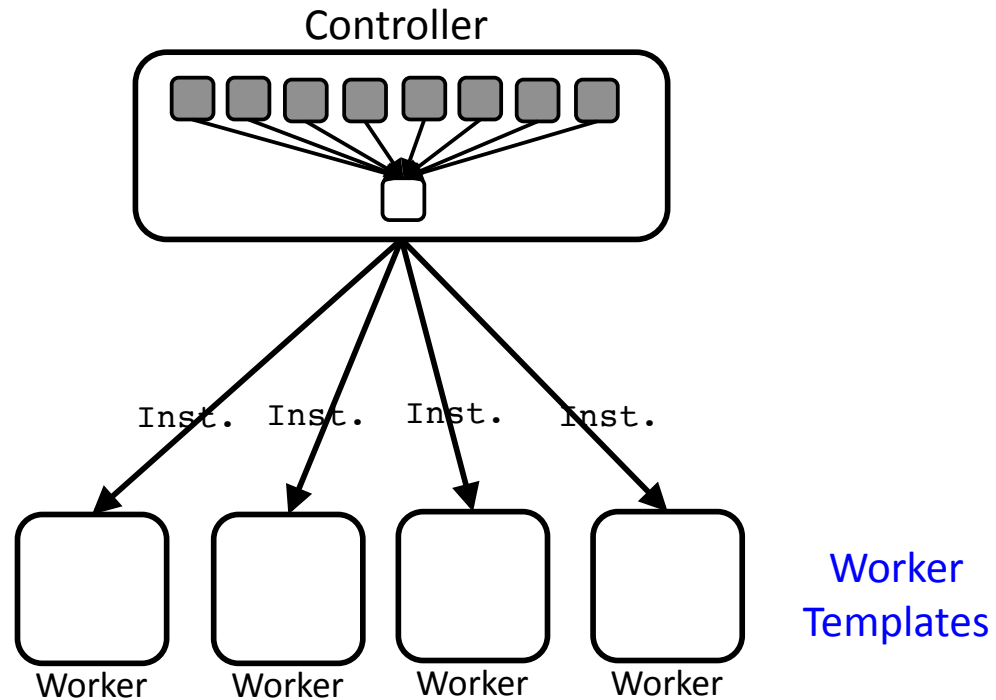
# Nimbus

## Controller and Worker Templates



# Nimbus

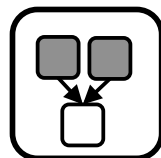
## Controller and Worker Templates



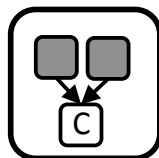
# Nimbus

## Controller and Worker Templates

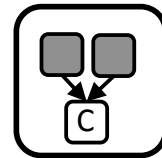
Controller



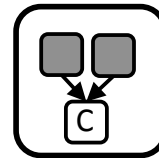
Worker



Worker



Worker

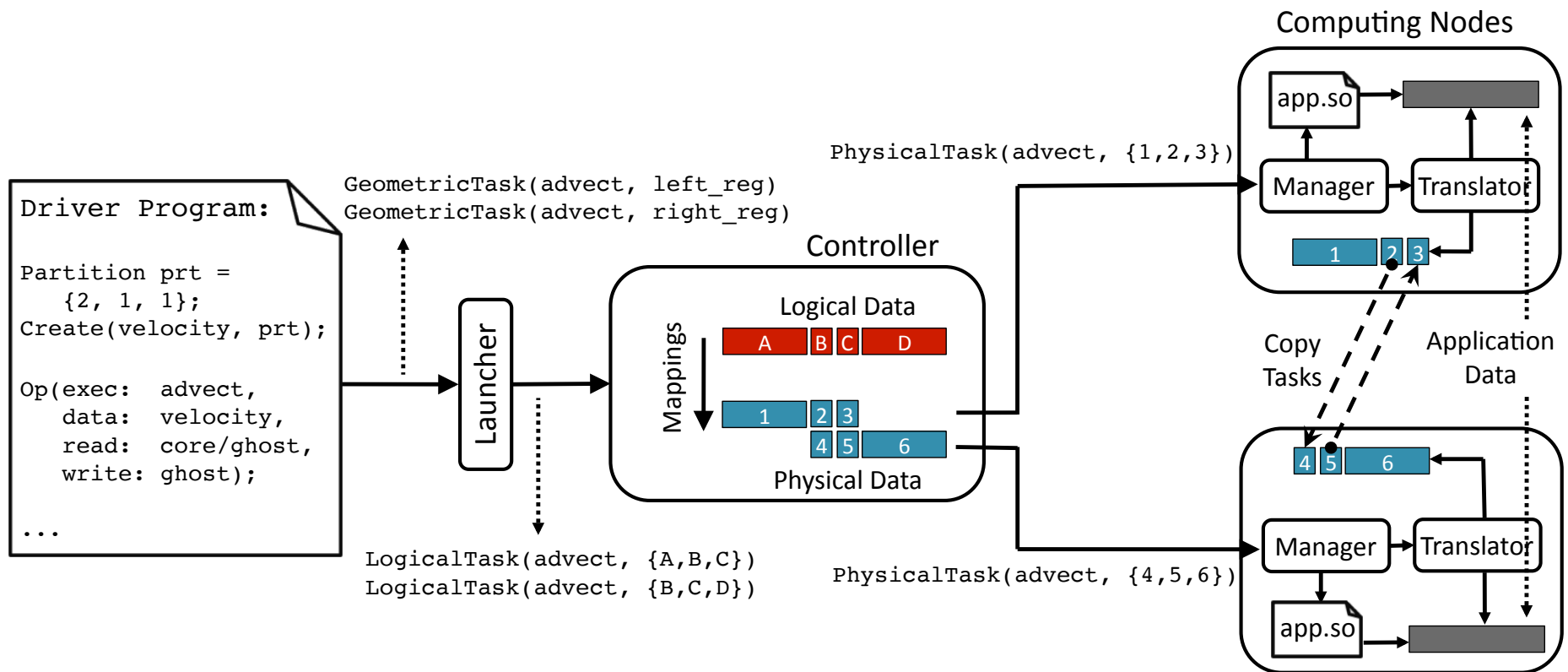


Worker

Worker  
Templates

# Nimbus

## Graphical Simulations



- The goal is to automatically distribute sequential library kernels.
- Four layer data abstraction (geometric, logical, physical, application).
- Automatic translation and caching between the data layers.



[nimbus.stanford.edu](https://nimbus.stanford.edu)

- For more information you can visit Nimbus website.



# This talk

- Control Plane: the Emerging Bottleneck
- Design Scope of the Control Plane
- Execution Templates
- Nimbus: a Framework with Templates
- Evaluation

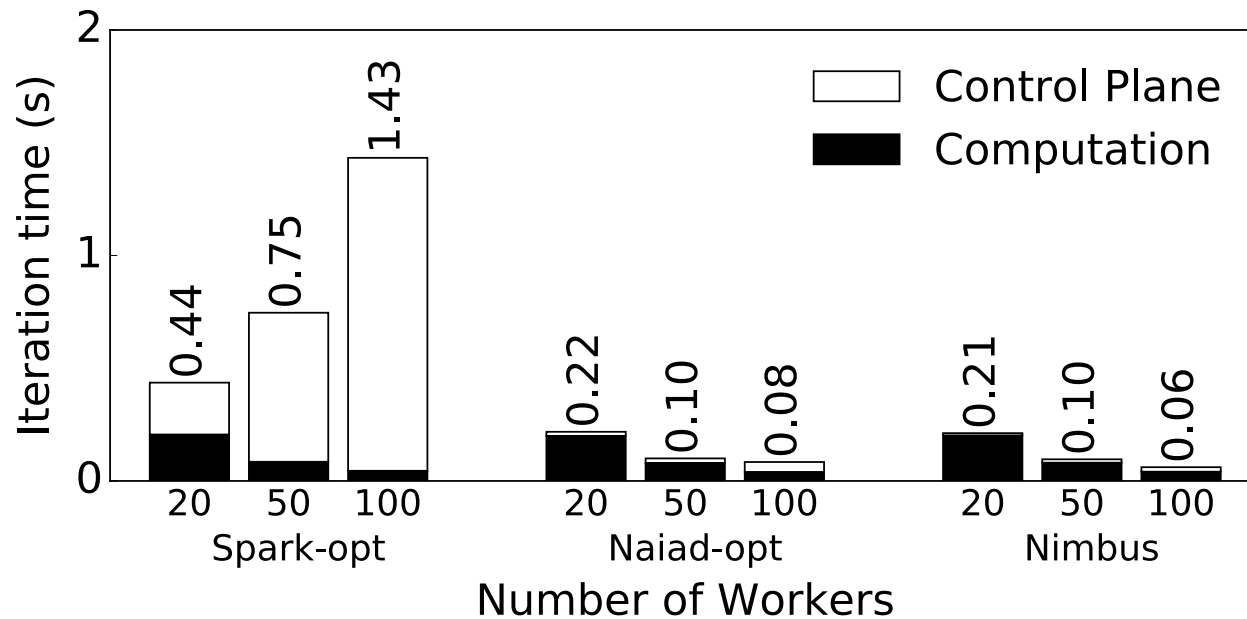
# Evaluation

## Results Summary

- Control plane task throughput:
  - Execution templates match the strong scaling performance of frameworks with distributed control plane design.
- Dynamic scheduling:
  - Execution templates allows low cost, reactive scheduling and dynamic resource allocation similar to a centralized frameworks.
- Dynamic control flow:
  - Execution templates can handle applications with nested loops and data dependent branches with low overhead.

# Evaluation

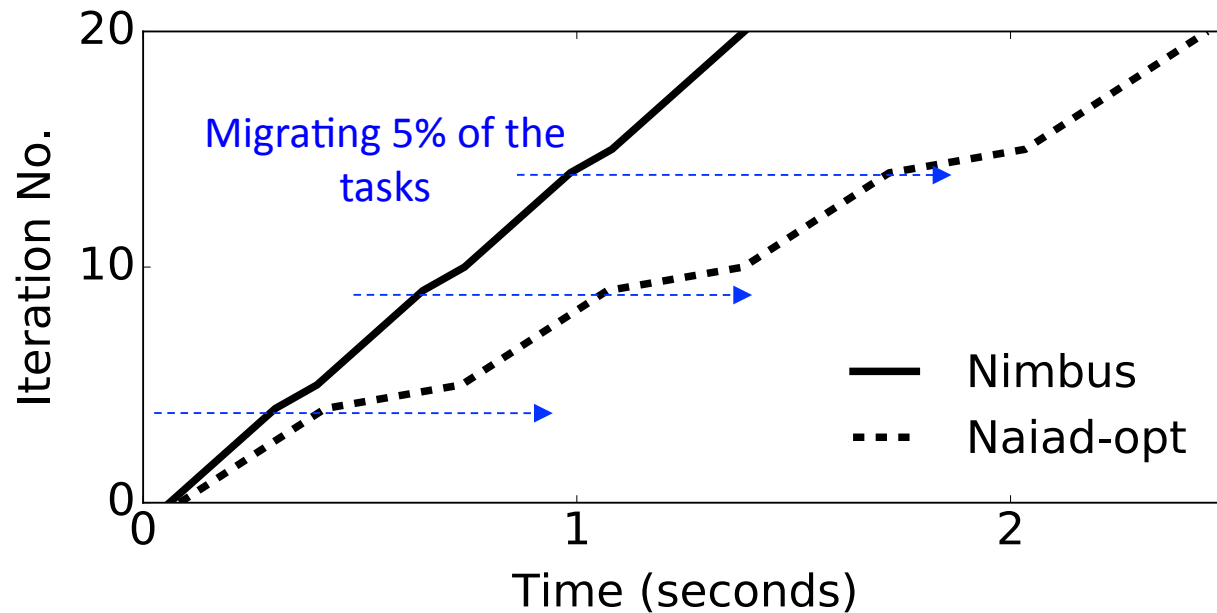
## Strong Scalability with Templates



- Logistic regression over data set of size 100GB.
- Spark-opt and Naiad-opt, runs tasks as fast as C++ implementation.
- Nimbus centralized controller with execution templates matches the performance of Naiad with a distributed control plane.

# Evaluation

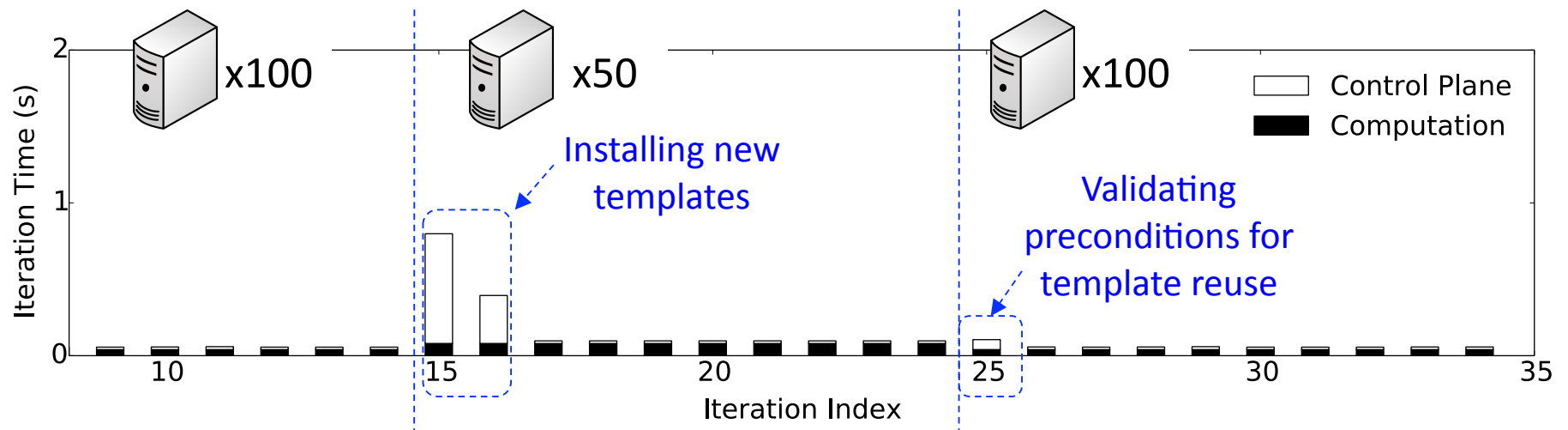
## Reactive, Fine-Grained Scheduling with Templates



- Logistic regression over data set of size 100GB, on 100 workers.
- Naiad-opt curve is simulated (migrations every 5 iterations).
- Execution templates allow low cost, reactive scheduling changes through edits at task granularity.
  - Single edit overhead is only  $41\mu\text{s}$  (in average).

# Evaluation

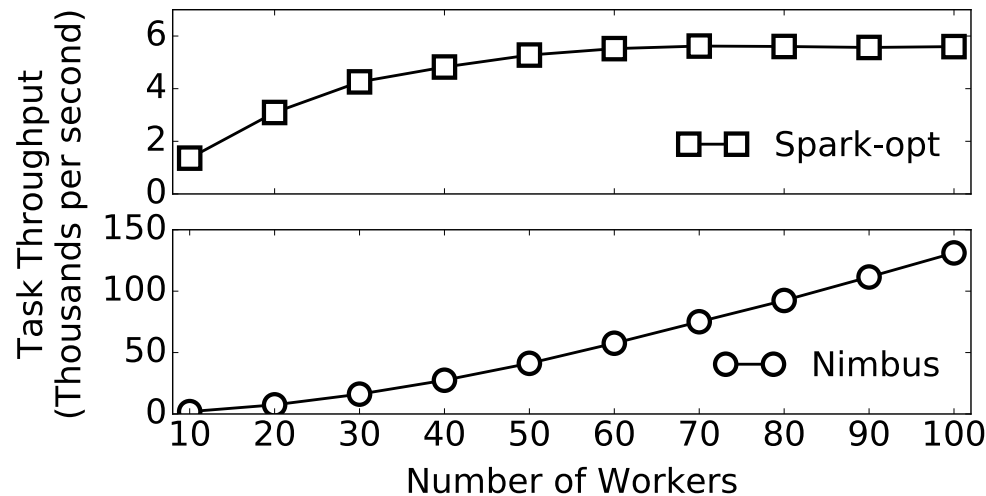
## Dynamic Resource Allocation with Templates



- Logistic regression over 100GB of data, on 50/100 workers.
- One-time template installation cost is ~40% of direct task scheduling.
- Nimbus allows dynamic resource allocation.
- Nimbus installs multiple versions of a template depending on resources.

# Evaluation

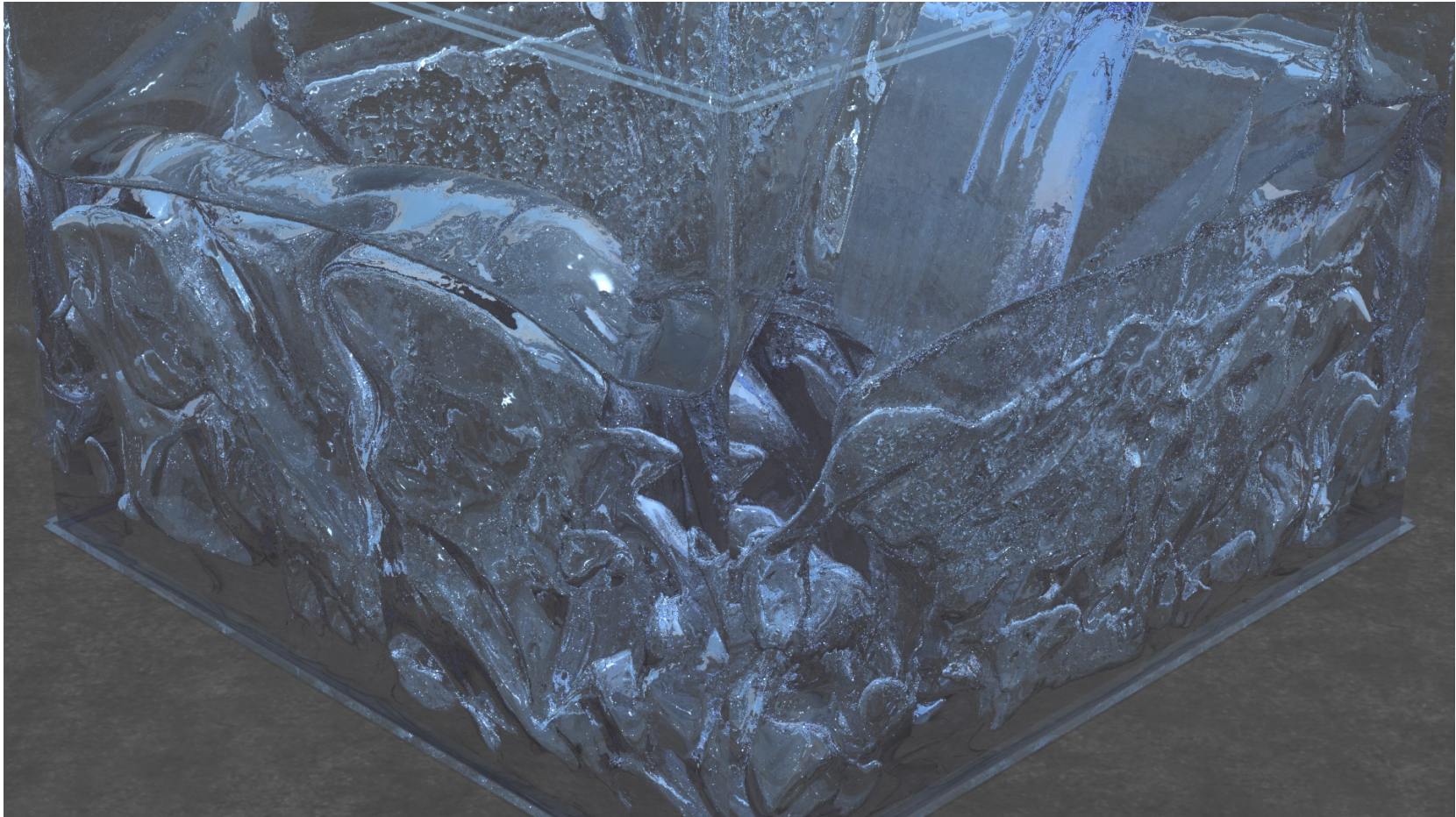
## High Task Throughput with Templates



- Spark and Nimbus both have centralized controller.
- Nimbus task throughput scales super linearly with more workers.
  - $O(N^2)$ : more tasks and shorter tasks, simultaneously.
- For a task graphs with single stage:
  - Instantiation cost is  $<2\mu s$  per task (500,000 tasks per second).

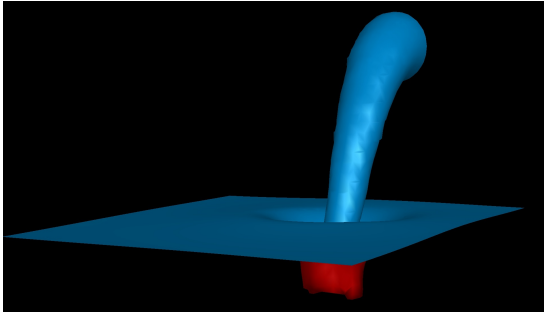
# Evaluation

Graphical Simulations Distributed in Nimbus

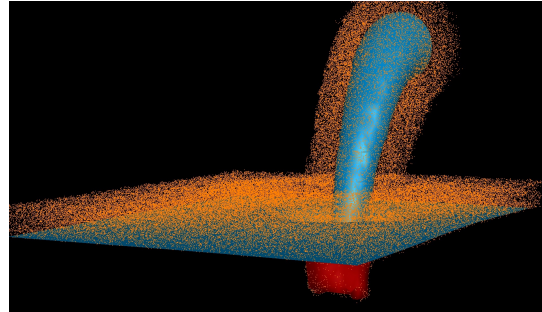


# Evaluation

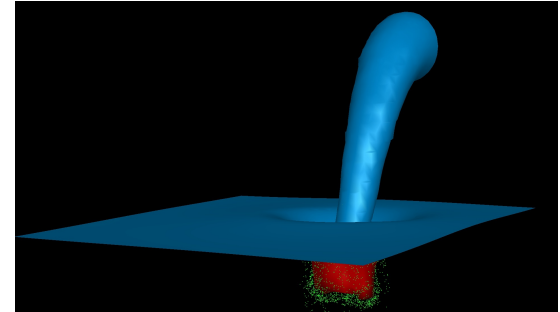
## Complexities of Graphical Simulations



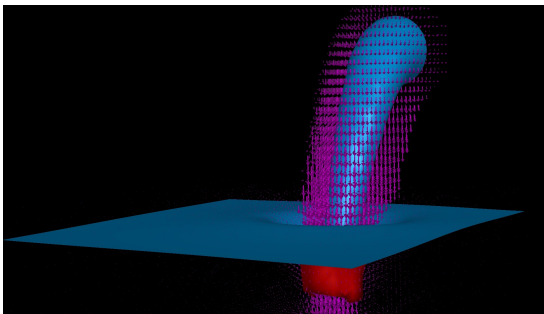
Levelset



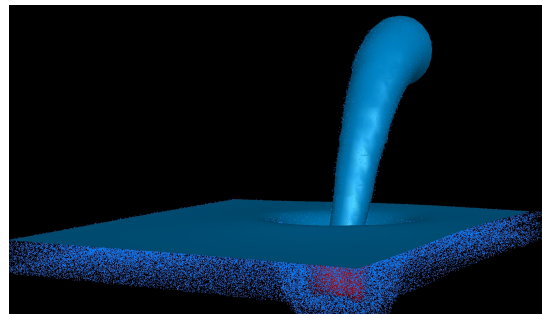
Positive Particles



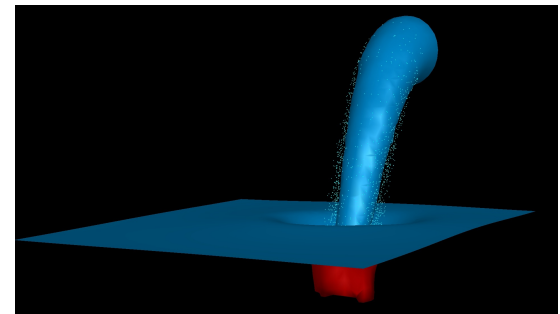
Positive Removed Particles



Velocity



Negative Particles



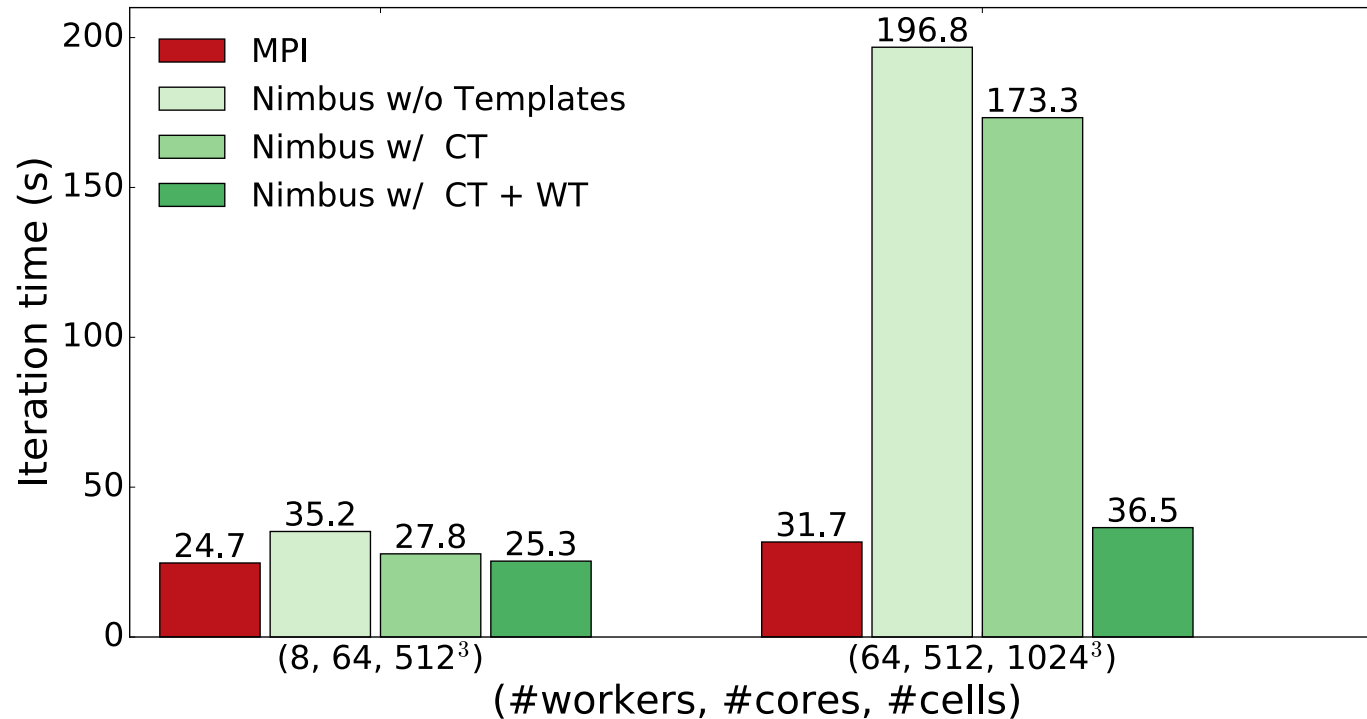
Negative Removed Particles

- 40 different variables: scalar, vector, particle.
- Triply nested loop with data dependent branches.
  - 9 different templates (basic blocks).
  - 3 branches that need patching.



# Evaluation

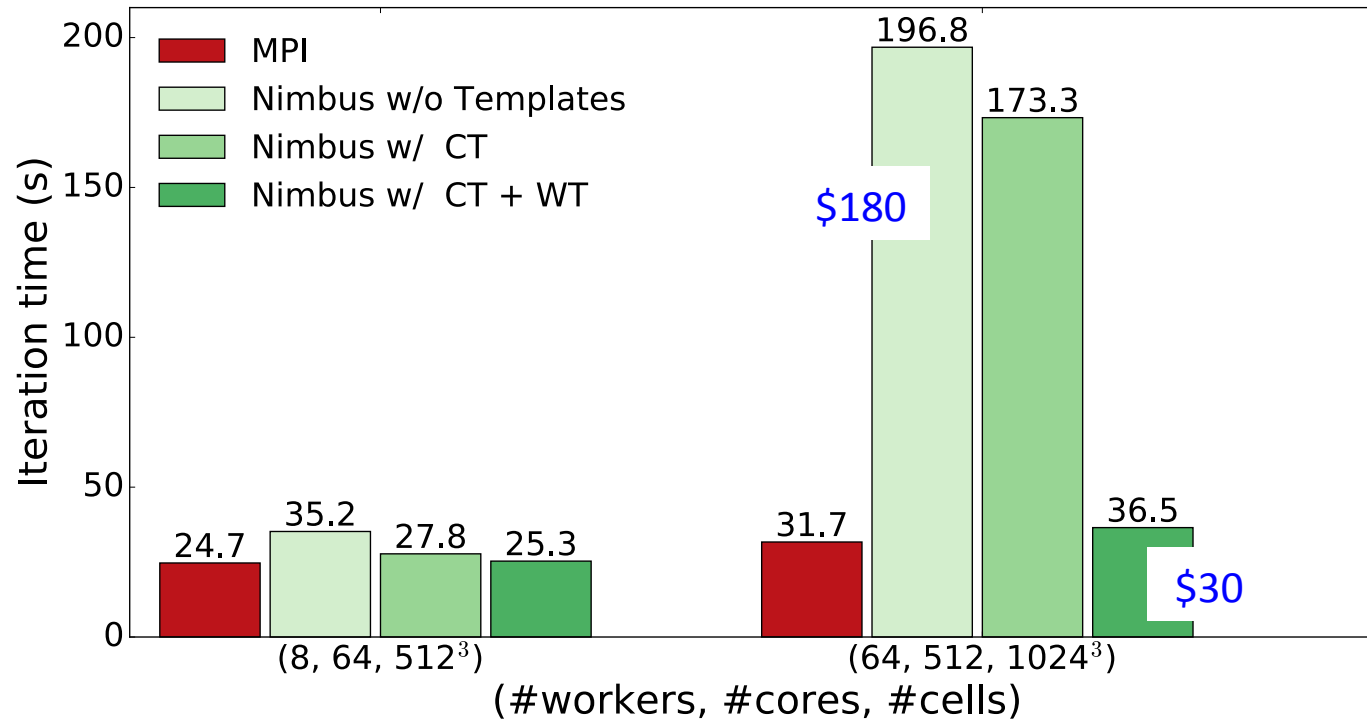
## Speedup with Templates



- Canonical water simulations under Nimbus and MPI.
- Without templates, Nimbus is almost 6x slower than MPI.
- Slow down means either lower resolution or more time/money.

# Evaluation

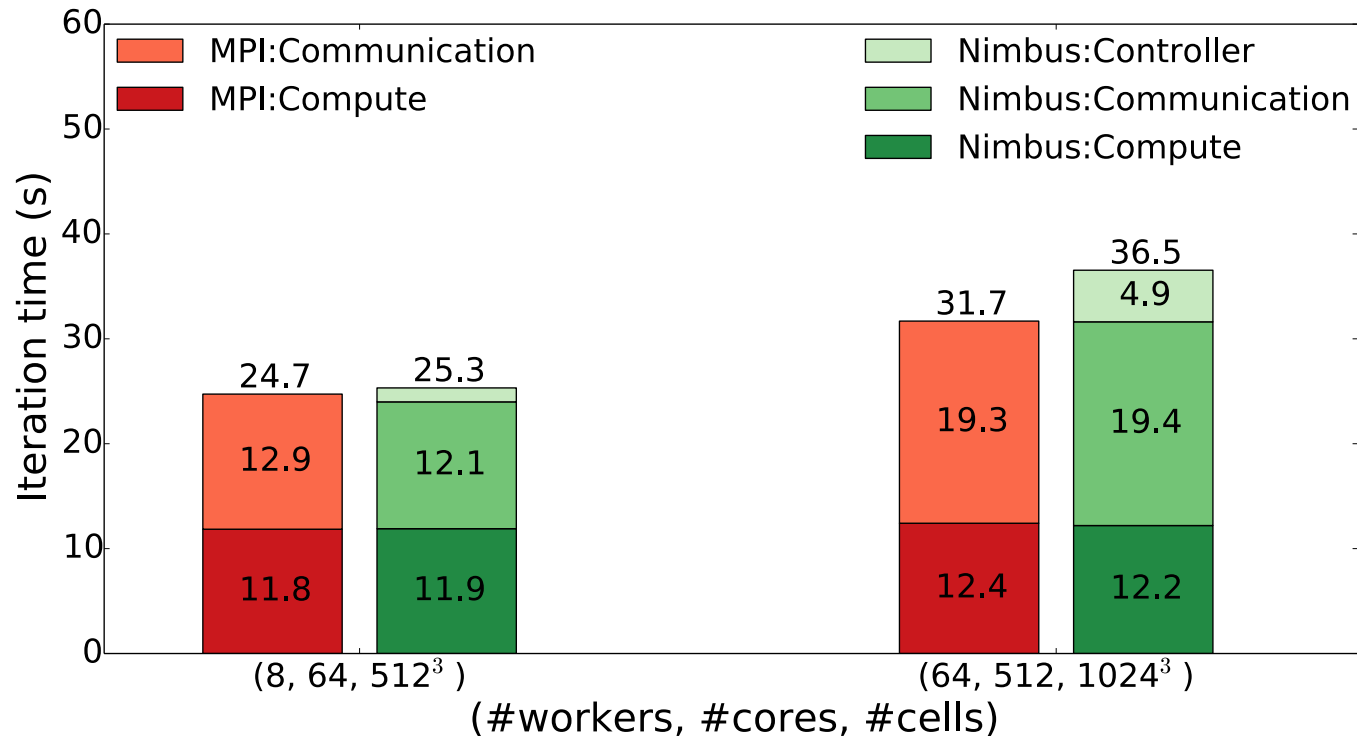
## Speedup with Templates



- Canonical water simulations under Nimbus and MPI.
- Without templates, Nimbus is almost 6x slower than MPI.
- Slow down means either lower resolution or more time/money.

# Evaluation

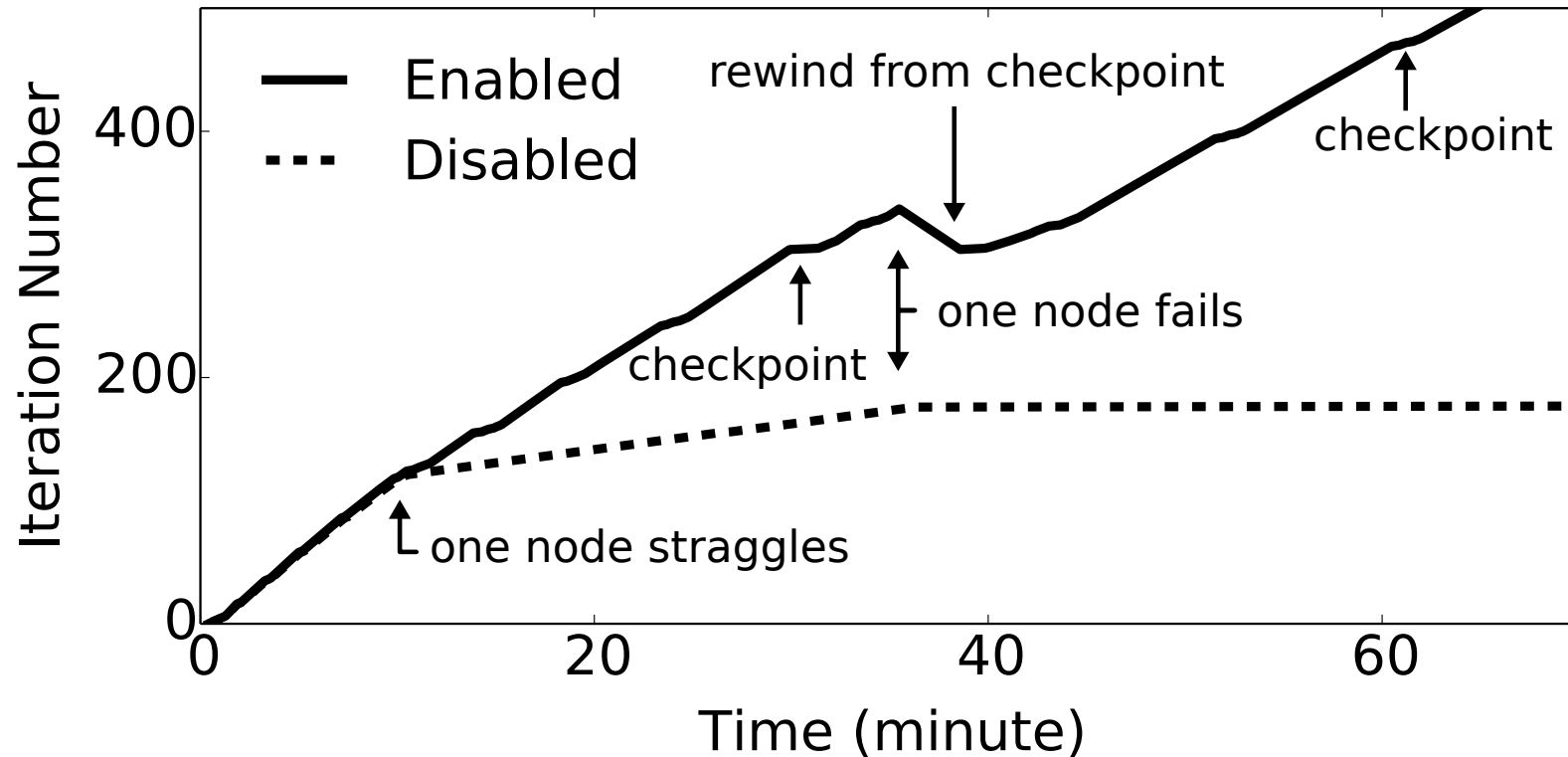
## Comparison with Hand-Tuned MPI



- Canonical water simulations under Nimbus and MPI.
- Nimbus performance is within 3-15% of the hand-tuned MPI.
- At 512 cores, there are more than **1 million** distinct data objects and task throughput picks at **460,000** tasks per second.

# Evaluation

## Load Balancing and Fault Recovery with Templates



- Nimbus controller adapts to the stragglers and worker failures.
- Templates are seamlessly installed as schedule changes.

# Contributions

- Demonstrating how the **control plane** is the emerging **bottleneck** for data analytics frameworks.
- **Execution Templates** as an abstraction for the control plane of cloud computing frameworks, that enables orders of magnitude higher task throughput, while keeping the fine-grained, flexible scheduling.
- The design, implementation, and evaluation of **Nimbus**, a distributed cloud computing framework that embeds execution templates.
- A demonstration of a single-core **graphical simulation** that Nimbus **automatically distributes** in the cloud showing execution templates in practice for complex applications.

# Conclusion

Control Plane Design	Example Framework	Task Throughput	Task Scheduling
Centralized	MapReduce	Low	Dynamic
	Hadoop		
	Spark		
Distributed	Naiad TensorFlow	High	Static
Centralized w/ Execution Templates	Nimbus	High	Dynamic

Thank You!