

Verified Progress Tracking for Timely Dataflow Matthias Brun, Sára Decova, Andrea Lattuada Department of Computer Science, ETH Zürich Dmitriv Traytel Department of Computer Science, University of Copenhagen

ETHzürich





Verified Progress Tracking for Timely Dataflow Matthias Brun, Sára Decova, Andrea Lattuada Department of Computer Science, ETH Zürich Dmitriv Traytel Department of Computer Science, University of Copenhagen

ETHzürich

A dataflow program



A dataflow program



Time-aware dataflow system

Timely Dataflow

Each tuple (message) tagged with timestamp

 logical/temporal grouping

Dataflow programming for automatic parallelzation

logical dataflow



Dataflow programming for automatic parallelzation

logical dataflow



physical dataflow









first propagation



first propagation





first propagation















iterative graph algorithm



add edge











































Progress Tracking

Core coordination component of Timely Dataflow

stream of label updates

Tracks outstanding work in the (i,0)(n,1)(m,3)(m,4)

Informs operators when all data with a certain timestamp has been delivered







extend the existing model (Abadi et al.) of the distributed component



extend the existing model (Abadi et al.) of the distributed component model, and verify, the worker-local propagation component



extend the existing model (Abadi et al.) of the distributed component model, and verify, the worker-local propagation component verify safety of the full Progress Tracker



Progress tracking concepts (1)

Progress tracking concepts (1)

(l, t) a datum at rest at an operator, or in motion on a channel Pointstamp





a datum at rest at an operator, or in motion on a channel

Progress tracking concepts (1) Pointstamp (l, t) a datum at rest at an oper

Timestamp

Location

a datum at rest at an operator, or in motion on a channel

semantic (causal) grouping of the data tupes of positive integers, or any type with ≤

Progress tracking concepts (1)

Pointstamp

a datum at rest at an operator, or in motion on a channel (l, t)Timestamp semantic (causal) grouping of the data tupes of positive integers, or any type with \leq

operator input or output port Location




Progress Tracking tracks the pointstamps and summarizes them as a Frontier represented as antichain F

a datum at rest at an operator, or in motion on a channel

semantic (causal) grouping of the data tupes of positive integers, or any type with ≤

operator input or output port



the operator may receive any timestamp t for which $\exists t' \in F$. $t' \leq t$

























by the operator



by the operator

Progress tracking

computes frontiers in two steps





the dataflow graph



distributed component **exchanges** pointstamp changes approximate, conservative view of all pointstamps

- worker-local component propagates changes on
 - update frontiers at operator input ports



Abadi et al.'s exchange protocol

Formal Analysis of a Distributed Algorithm for Tracking Progress

Martín Abadi^{1,2}, Frank McSherry¹, Derek G. Murray¹, and Thomas L. Rodeheffer¹

¹ Microsoft Research Silicon Valley
 ² University of California, Santa Cruz

Abstract. Tracking the progress of computations can be both important and delicate in distributed systems. In a recent distributed algorithm for this purpose, each processor maintains a delayed view of the pending work, which is represented in terms of points in virtual time. This paper presents a formal specification of that algorithm in the temporal logic TLA, and describes a mechanically verified correctness proof of its main properties.

- Formalized in TLA+PS before
- We ported the formalization to Isabelle
- One person-week, Sketch & Hammer



Abadi et al.'s exchange protocol

Formal Analysis of a Distributed Algorithm for **Tracking Progress**

Martín Abadi^{1,2}, Frank McSherry¹, Derek G. Murray¹, and Thomas L. Rodeheffer¹

> ¹ Microsoft Research Silicon Valley ² University of California, Santa Cruz

Abstract. Tracking the progress of computations can be both important and delicate in distributed systems. In a recent distributed algorithm for this purpose, each processor maintains a delayed view of the pending work, which is represented in terms of points in virtual time. This paper presents a formal specification of that algorithm in the temporal logic TLA, and describes a mechanically verified correctness proof of its main properties.

- Formalized in TLA+PS before
- We ported the formalization to Isabelle
- One person-week, Sketch & Hammer
- **Discovered problems**:
- 1. Transitions based on global state
- 2. Too strong restrictions on allowed transitions





Abadi et al.'s exchange protocol

Formal Analysis of a Distributed Algorithm for **Tracking Progress**

Martín Abadi^{1,2}, Frank McSherry¹, Derek G. Murray¹, and Thomas L. Rodeheffer¹

> ¹ Microsoft Research Silicon Valley ² University of California, Santa Cruz

Abstract. Tracking the progress of computations can be both important and delicate in distributed systems. In a recent distributed algorithm for this purpose, each processor maintains a delayed view of the pending work, which is represented in terms of points in virtual time. This paper presents a formal specification of that algorithm in the temporal logic TLA, and describes a mechanically verified correctness proof of its main properties.

- Formalized in TLA+PS before
- We ported the formalization to Isabelle
- One person-week, Sketch & Hammer
- **Discovered problems**:
- 1. Transitions based on global state
- 2. Too strong restrictions on allowed transitions





rec :: 'p zmset temp :: $'w \Rightarrow 'p \ zmset$ glob :: ' $w \Rightarrow$ 'p zmset

```
record ('w :: finite, 'p :: order) conf =
msg :: 'w \Rightarrow 'w \Rightarrow 'p \ zmset \ list
```



Finite set of workers rec :: 'p zmset temp :: $'w \Rightarrow 'p \ zmset$ glob :: $'w \Rightarrow 'p \ zmset$

```
record ('w :: finite, 'p :: order) conf =
msg :: 'w \Rightarrow 'w \Rightarrow 'p \ zmset \ list
```



Finite set of workers **record** ('w :: finite, 'p :: order) conf =rec :: 'p zmset $msg :: 'w \Rightarrow 'w \Rightarrow 'p \ zmset \ list$ temp :: $'w \Rightarrow 'p \ zmset$ glob :: ' $w \Rightarrow$ 'p zmset

Partially-ordered pointstamps



Finite set of workers \mathbf{record} ('w :: finiterec :: 'p zmse $msg :: 'w \Rightarrow$ temp :: $'w \Rightarrow 'p \ zmset$ glob :: ' $w \Rightarrow$ 'p zmset

$$et < signed multiset$$

 $w \Rightarrow 'p \ zmset \ list$



$$et < signed multiset$$

 $w \Rightarrow 'p \ zmset \ list$



$$et < signed multiset$$

 $w \Rightarrow 'p \ zmset \ list$





 $c' = c (\operatorname{rec} = \operatorname{rec} c + \Delta, \operatorname{temp} = (\operatorname{temp} c)(w := \operatorname{temp} c w + \Delta))$

Partially-ordered
pointstampsnite, 'p :: order) conf =nite, 'p :: order) conf =isigned multiset'w
$$\Rightarrow$$
 'p zmset list'p zmset'p zmsetNot-yet-communicated
pointstamp updatesch worker's local
ation of the global state

definition perf_op :: ' $w \Rightarrow$ 'p mset \Rightarrow 'p mset \Rightarrow ('w, 'p) conf \Rightarrow ('w, 'p) conf \Rightarrow bool where perf_op $w \Delta_{neg} \Delta_{pos} c c' = \underline{\text{let}} \Delta = \Delta_{pos} - \Delta_{neg} \underline{\text{in}} (\forall p. \Delta_{neg} \# p \leq \text{rec} c \#_z p) \land \text{upright} \Delta \land$

 $c' = c (\operatorname{rec} = \operatorname{rec} c + \Delta, \operatorname{temp} = (\operatorname{temp} c)(w := \operatorname{temp} c w + \Delta))$

Partially-ordered
pointstampsnite, 'p :: order) conf =nite, 'p :: order) conf =signed multiset'w
$$\Rightarrow$$
 'p zmset list'p zmset'p zmsetNot-yet-communicated
pointstamp updatesch worker's local
ation of the global state

definition perf_op :: ' $w \Rightarrow 'p \ mset \Rightarrow 'p \ mset \Rightarrow ('w, 'p) \ conf \Rightarrow ('w, 'p) \ conf \Rightarrow bool where$ perf_op $w \Delta_{neg} \Delta_{pos} c c' = \underline{let} \Delta = \Delta_{pos} - \Delta_{neg} \underline{in} (\forall p. \Delta_{neg} \# p \leq \underline{rec} c \#_z p) \land upright \Delta \land$

Addresses both problems:

- 2. Allow creation of pointstamps based on capabilities

Our exchange protocol

1. Split global **rec** into worker-local capabilities caps :: $'w \Rightarrow 'p \ zmset$

Addresses both problems:

- 1. Split global **rec** into worker-local capabilities caps :: $'w \Rightarrow 'p \ zmset$
- 2. Allow creation of pointstamps based on capabilities

We verify the **same** safety property as Abadi et al.

Our exchange protocol

- If any worker's **glob** becomes vacant up to some pointstamp, then that pointstamp and any lesser ones do not exist in the system, i.e., are not present in **rec** (and will remain so).

Worker-local propagation

Dataflow \approx directed graph with antichains (over a partial order) as weights A well-behaved operation combining timestamps and summaries

Worker-local propagation

Safety property:

If all worklists neither contain timestamp t (adjusted by summaries) nor smaller timestamps, then all locations know whether they may encounter t (adjusted by summaries) in the future.

Combined progress tracker

19


Not what Timely Dataflow does. Rather:



Termination of propagation not verified



Not what Timely Dataflow does. Rather:



Termination of propagation not verified

Combined safety property:

Every initialized worker w has some evidence for the existence of a timestamp t at location l at any worker w' in w's frontier at all locations l' reachable from l.

Summary

Modeled two components



Formalized & verified safety properties for both

Verified safety of the combined Progress Tracker

https://www.isa-afp.org/entries/Progress_Tracking.htm

Í.
1

r	Υ	וו

\mathbf{O}	\bigcap
2	U

Summary

Modeled two components



Formalized & verified safety properties for both

Verified safety of the combined Progress Tracker

https://www.isa-afp.org/entries/Progress_Tracking.html//www.isa-Afp.org/entries/Progress_Tracking.html//www.isa-Afp.org/entrie



L.	
	1
-	

r	Υ	ר



Summary

Modeled two components



Formalized & verified safety properties for both

Verified safety of the combined Progress Tracke

https://www.isa-afp.org/entries/Progress_Tracking.html



	Pr	oof as	ssistar	nt tec	hnolo
		Local type cla	es & asses	Sledgehammer	
er			Signed M	lultisets	abelle
ml					Nº CR
)					
У					





Modeled two components



Formalized & verified safety properties for both





Verified Progress Tracking for Timely Dataflow Matthias Brun, Sára Decova, Andrea Lattuada Department of Computer Science, ETH Zürich Dmitriv Traytel Department of Computer Science, University of Copenhagen

ETHzürich



Verified Progress Tracking for Timely Dataflow Matthias Brun, Sára Decova, Andrea Lattuada Department of Computer Science, ETH Zürich Dmitriy Traytel Thank you! Department of Computer Science, University of Copenhagen **Questions?**

ETHzürich