The Semantics of Weak Persistency

Viktor Vafeiadis



CALCO, 3 September 2021



ERC-COG-2020 grant PERSIST, Mar'21-Feb'26

Traditional computer storage



RAM -RANDOM ACCESS MEMORY (MEMORY)

fast byte-size access volatile storage



HDD -HARD DISK DRIVE (STORAGE)

- slow block-size access
- durable storage

Non-volatile memory (NVM)

Fast byte-size access

- ▶ 3-10x slower than RAM
- 100x faster that HDD

Durable storage

- Larger than RAM
- More energy-efficient than RAM















Roadmap

- What semantics do programs have? (operational, declarative/axiomatic)
- How can we avoid weak behaviours? (flushes, fences)
- When is a persistent algorithm correct? (invariants, persistent serializability/linearizability)
- What techniques can we use to establish correctness? (program logics, model checking)

Weak persistency semantics



40s Sequential programs

40s Sequential programs

70s

Interleaving concurrency (SC)



Weak memory consistency

The x86-TSO model



Store buffering (SB) Initially, x = y = 0. x := 1; $a := y \parallel 0$ y := 1; $b := x \parallel 0$

40s	Sequential programs
70s	Interleaving concurrency (SC)
90s	Weak memory concurrency (WMC)
now	WMC & Weak memory persistency





!! Execution continues *ahead of persistence*

- asynchronous persists



!! Execution continues ahead of persistence

- asynchronous persists

Writes may persist out of order — relaxed persists

Consistency Model

the **order** in which writes are **made visible** to other threads

Persistency Model

the *order* in which writes are *persisted* to NVM

NVM Semantics Consistency + Persistency Model



Basic persistency model



x := 1 : add x := 1 to p-buffer

a:=x : if p-buffer contains x, read latest entry else read from memory

p-buffer lost; memory retained

Unbuffered at *non-deterministic* points in time

Buffering & unbuffering orders may disagree

Relaxed persists



!! out of order persists

Explicit persists?



!! out of order persists

explicit persists?

x86: clwb, clflushopt, clflush



- * clwb and clflushopt: same ordering constraints
- * clwb does not invalidate cache line
- * clflushopt invalidates cache line
- * clflush: strongest ordering constraints; invalidates cache line

Strong persists: clflush



Weak persists: clflushopt & clwb



Weak persists: clflushopt & clwb



Solution: Persist sequences



Adding weak memory consistency







buffer/unbuffer order: consistency model





Two equivalent formal models

and

Operational



Declarative/Axiomatic



Roadmap

- What semantics do programs have? (operational, declarative/axiomatic)
- How can we avoid weak behaviours? (flushes, fences)
- When is a persistent algorithm correct? (invariants, persistent serializability/linearizability)
- What techniques can we use to establish correctness? (program logics, model checking)

Serialisability (SER)

All transactions appear to execute in a sequential order



All transactions appear to execute in a sequential order

A prefix of transactions appears to persist in the same sequential order



All transactions appear to execute in a sequential order

A prefix of transactions appears to persist in the same sequential order



All transactions appear to execute in a sequential order

A prefix of transactions appears to persist in the same sequential order



All transactions appear to execute in a sequential order

A prefix of transactions appears to persist in the same sequential order

in <u>each era</u>



All transactions appear to execute in a sequential order

A prefix of transactions appears to **persist** in the **same sequential order** in **each era**



All transactions appear to execute in a sequential order

A prefix of transactions appears to persist in the same sequential order

in **each era**



Is PSER *Feasible*?

✓ PSER *implementation* in *x86, ARM*

Take SER Implementation — e.g. 2-PL

- + add code for *persistence*
- + add code to *log metadata* for *recovery*
- + add recovery mechanism



Is PSER Useful?

Given library *L* (e.g. queue library):

1. Take any correct sequential implementation of L

2. wrap each operation in a PSER transaction

 \Rightarrow correct, concurrent & persistent implementation of L



sequential queue imp.

enq(q,v)=	
pser{ <enq_body></enq_body>	}
deq(q)= pser{	
<deq_body></deq_body>	}

correct concurrent & persistent queue imp.

Roadmap

- What semantics do programs have? (operational, declarative/axiomatic)
- How can we avoid weak behaviours? (flushes, fences)
- When is a persistent algorithm correct?
 (invariants, persistent serializability/linearizability)
- What techniques can we use to establish correctness? (program logics, model checking)

Program logics for persistency

Instrument program with additional variables:

- volatile x_v the latest observable value of x
- persisted x_p the persisted value of x

Use OGRA, a slight weakening of Owicki-Gries that is sound under x86-TSO



See OOPSLA'20; the presentation here is somewhat simplified

Instrument program with additional variables:

- \blacktriangleright volatile x_{y} the latest observable value of x
- \blacktriangleright persisted $x_{\rm p}$ the persisted value of x

Use OGRA, a slight weakening of Owicki-Gries that is sound under x86-TSO



See OOPSLA'20: the presentation here is somewhat simplified

 $X_p := X_v$

 $y_n := y_n$

Instrument program with additional variables:

- volatile x_v the latest observable value of x
- persisted x_p the persisted value of x

Use OGRA, a slight weakening of Owicki-Gries that is sound under x86-TSO

 $egin{aligned} x_{
u} &:= 1 \ x_{
u} &:= x_{
u} \ y_{
u} &:= 1 \end{aligned}$

See OOPSLA'20; the presentation here is somewhat simplified

Parallel context:

 $X_p := X_v$

 $y_n := y_n$

Instrument program with additional variables:

- volatile x_v the latest observable value of x
- persisted x_p the persisted value of x

Use OGRA, a slight weakening of Owicki-Gries that is sound under x86-TSO

$$egin{aligned} & \{x_{
u}=0 \wedge x_{
u}=0 \wedge y_{
u}=0 \wedge y_{
u}=0 \} \ & x_{
u}:=1 \ & x_{
u}:=x_{
u} \ & y_{
u}:=1 \end{aligned}$$

Parallel context:

$$\begin{array}{l} x_{p} := x_{v} \\ y_{p} := y_{v} \end{array}$$

See OOPSLA'20; the presentation here is somewhat simplified

Instrument program with additional variables:

- volatile x_v the latest observable value of x
- persisted x_p the persisted value of x

Use OGRA, a slight weakening of Owicki-Gries that is sound under x86-TSO

$$\{ x_{\nu} = 0 \land x_{p} = 0 \land y_{\nu} = 0 \land y_{p} = 0 \}$$

$$x_{\nu} := 1$$

$$\{ x_{\nu} = 1 \land x_{p} \in \{0, 1\} \land y_{\nu} = 0 \land y_{p} = 0 \}$$

$$x_{p} := x_{\nu}$$

$$y_{\nu} := 1$$

See OOPSLA'20; the presentation here is somewhat simplified

Parallel context:

 $\begin{array}{l} x_p := x_v \\ y_p := y_v \end{array}$

Instrument program with additional variables:

- volatile x_v the latest observable value of x
- persisted x_p the persisted value of x

Use OGRA, a slight weakening of Owicki-Gries that is sound under x86-TSO

$$\{ x_{\nu} = 0 \land x_{p} = 0 \land y_{\nu} = 0 \land y_{p} = 0 \}$$

$$x_{\nu} := 1$$

$$\{ x_{\nu} = 1 \land x_{p} \in \{0, 1\} \land y_{\nu} = 0 \land y_{p} = 0 \}$$

$$x_{\rho} := x_{\nu}$$

$$\{ x_{\nu} = 1 \land x_{p} = 1 \land y_{\nu} = 0 \land y_{p} = 0 \}$$

$$y_{\nu} := 1$$

Parallel context:

$$\begin{array}{l} x_{p} := x_{v} \\ y_{p} := y_{v} \end{array}$$

See OOPSLA'20; the presentation here is somewhat simplified

Instrument program with additional variables:

- volatile x_v the latest observable value of x
- persisted x_p the persisted value of x

Use OGRA, a slight weakening of Owicki-Gries that is sound under x86-TSO

$$\{ x_{\nu} = 0 \land x_{p} = 0 \land y_{\nu} = 0 \land y_{p} = 0 \}$$

$$x_{\nu} := 1$$

$$\{ x_{\nu} = 1 \land x_{p} \in \{0, 1\} \land y_{\nu} = 0 \land y_{p} = 0 \}$$

$$x_{\rho} := x_{\nu}$$

$$\{ x_{\nu} = 1 \land x_{p} = 1 \land y_{\nu} = 0 \land y_{p} = 0 \}$$

$$y_{\nu} := 1$$

$$\{ x_{\nu} = 1 \land x_{p} = 1 \land y_{\nu} = 1 \land y_{p} \in \{0, 1\} \}$$

Parallel context:

$$\begin{array}{l} x_{p} := x_{v} \\ y_{p} := y_{v} \end{array}$$

See OOPSLA'20; the presentation here is somewhat simplified

Model checking for persistency

Software model checking

Given a program P and a property Φ : Check that all executions of P satisfy Φ .

- P shared-memory concurrency + flushes/fences
- \blacktriangleright Φ safety assertion, invariant over persisted state
- Stateless' enumeration of all reachable program states

Which kind of model to use?



or



Which kind of model to use?

or



Declarative [init] Wx1 Wy1 Wz1Wz1 Wy1 Wz1

Rv1

 R_z1

 $R \times 1$

- Easy to build an interpreter
- X Tied to a specific model
- X Many redundant executions

- X Not immediately executable
- ✓ Largely model-agnostic
- ✓ Almost no redundancy

Which kind of model to use?

Operational CDII CDU Initially, x = y = z = 0x := 1 || y := 1 || z := 1 Or a := x || b := y || c := znersist Has $6!/2^3 = 90$ interleavings but only one execution graph X Tied to a specific model Many redundant executions

Declarative



× Not immediately executable

- ✓ Largely model-agnostic
- ✓ Almost no redundancy

Using declarative models pays off!

RMEM benchmarks			
	RMEM	GenMC	
DQ/211-2-1	172.34	0.17	
DQ-opt/211-2-1	770.45	0.17	
STC/210-011-000	1100.35	0.06	
STC-opt/210-011-000	1154.68	0.10	
QU/100-100-010	1099.20	0.06	
QU-opt/100-100-010	*	0.06	
Verification times in seconds			

- RMEM operational ARM model with a few optimisations
- GenMC declarative memory model

Construct all consistent execution graphs incrementally

Fix insertion order (e.g. increasing thread ID order)

Construct all consistent execution graphs incrementally

Fix insertion order (e.g. increasing thread ID order)

$$x := 1 \parallel a := x$$

Wx0

Construct all consistent execution graphs incrementally

Fix insertion order (e.g. increasing thread ID order)

$$x := 1 \parallel a := x$$



Construct all consistent execution graphs incrementally

Fix insertion order (e.g. increasing thread ID order)

When adding a read *r*:

Consider all possible writes that r could read from.

$$x := 1 \parallel a := x$$



Construct all consistent execution graphs incrementally

Fix insertion order (e.g. increasing thread ID order)

When adding a read *r*:

Consider all possible writes that r could read from.



Construct all consistent execution graphs incrementally

Fix insertion order (e.g. increasing thread ID order)

When adding a read *r*:

Consider all possible writes that r could read from.

Construct all consistent execution graphs incrementally

Fix insertion order (e.g. increasing thread ID order)

When adding a read *r*:

Consider all possible writes that r could read from.

When adding a write w:

Revisit existing reads to instead read from w.

$$\begin{array}{cccc} x := 1 & \| & a := x \end{array} & \text{Add } a := x \text{ first} \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & &$$

Construct all consistent execution graphs incrementally

Fix insertion order (e.g. increasing thread ID order)

When adding a read *r*:

Consider all possible writes that r could read from.

When adding a write w:

Revisit existing reads to instead read from w.



$$x := 1 \parallel y := 2 \parallel z := 3$$

Check that $y \neq 1$ after a crash

$$x := 1 \parallel y := 2 \parallel z := 3$$

Check that $y \neq 1$ after a crash

Simple approach

- Enumerate all post-crash states of a program
- Check that each persisted state satisfies the assertion

Total p-ordering

- In which order did the durable events persist?
- Take any prefix of that order
- ln total, $4 \times 3! = 24$ cases

$$x := 1 \parallel y := 2 \parallel z := 3$$

Check that $y \neq 1$ after a crash

Simple approach

- Enumerate all post-crash states of a program
- Check that each persisted state satisfies the assertion

Total p-ordering

- In which order did the durable events persist?
- Take any prefix of that order
- ln total, $4 \times 3! = 24$ cases

Partial p-ordering

- Which durable events have persisted?
- Their relative ordering is irrelevant
- ln total, $2^3 = 8$ cases

$$x := 1 \parallel y := 2 \parallel z := 3$$

Check that $y \neq 1$ after a crash

Simple approach

 \blacktriangleright In total. 4 \times 3! = 24 cases

- Enumerate all post-crash states of a program
- Check that each persisted state satisfies the assertion

Total p-orderin
In which order
Take any prefix of that order

ln total, $2^3 = 8$ cases

Recovery observer

Treat the persistency assertions as an additional thread

$$x := 1 \parallel y := 2 \parallel z := 3$$

 \checkmark $a := y$
 $assert(a \neq 1)$

... with somewhat different consistency axioms

Enforce snapshot atomicity (all recovery reads of y read from the same write)

 \blacktriangleright Observing a durable event \Rightarrow observe all events persisted before it

In our example, there are only 2 execution graphs

Summary

- What semantics do programs have? (operational, declarative/axiomatic)
- How can we avoid weak behaviours? (flushes, fences)
- When is a persistent algorithm correct?
 (invariants, persistent serializability/linearizability)
- ✓ What techniques can we use to establish correctness? (program logics, model checking)