# pBseq: probabilistic buffer-swapping sequence locks

Nicholas Mc Guire <safety@osadl.org>

October 21, 2016

# Outline

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

- Context - seqlocks and RT
- Alternatives
- pBseq
- Entropy harvesting
- Conclusions

# sequence locks overview

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

**Context**

Alternatives

pBseq

Inherent
randomness

Conclusions

```
single writer:          multiple readers:

                        unsigned seq_start, seq_end;
                        do {
                                seq_start = seq;
  data1 = ...;                  local1 = data1;
  data2 = ...;                  local2 = data2;
  seq = seq + 2;                seq_end = seq;
                        } while (seq_start != seq_end
                                 || seq_start % 2);
```

# sequence locks properties

- Writer never blocks on reader (no writer starvation)
- Reader may spin on writer (but no reader starvation)
- Multiple readers permitted
- Single writer (multi-writer -> by locks)
- Writers preferred over readers

Sequence locks have been in use since the late 2.5 kernel series,

# The seqlock problems for RT

- Reader retries are not bounded - even though generally short
- Anonymous locks do not permit boosting
- Readers/writers can't be boosted -> unbounded delays on writer preemption
- Long spinning readers -> cache impact

# Current mitigation

- Basic options
    - Try to add owner concept to sequence locks
    - Revert to boostable locks
- 1st option would be complex and unnecessary for non-RT
- 2nd is simple but limits scalability of RT
    - write_seqlock: grabs the seqlock spin_lock (multiple writer serialization)
    - read_seqbegin: spin_lock(&sl->lock);
      spin_unlock(&sl->lock); so it is bostable now.

No real mitigation in the current RT patch-set - the its more a workaround.

# Alternative mitigation: replication

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

Context

**Alternatives**

pBseq

Inherent
randomness

Conclusions

- SeqLocks are unbounded due to possible writer preemption
- Mitigation: replicated data
- Implementations:
    - seqcount_latch (fast ktime and in latch_tree_OP)
    - Suitable for non-atomic modifications
    - Prime motivation: unconditional lookups e.g. NMI context

For the NMI case data duplication is sufficient. This also could solve some of the RT issues but its not a simple replacement (code-level) for seqlocks.

# Concept of seqcount]_latch

- Maintain replicated data - data[0]/data[1]
- Redirect readers to stable copy
- Use LSB of sequence to select data buffer
- Non-probabilistic approach - redirection is deterministic
- Retry probability is assumed to be negligible due to preemption
- Worst-case: still unbounded - lockstep update/read possible

Lockstep behavior on larger multi-core is actually possible since cross-core delays can be quite large.

# pBseq probabilistic buffer-swapping sequence locks - overview

ⓈSADL

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

Context

Alternatives

pBseq

Inherent
randomness

Conclusions

# Concept of pBseq

- Maintain replicated data - data[NUM_REPLICA]
- Probabilistic approach - redirect readers to random copy
- Use inherent nondeterminism to select data copy
- Retry probability is statistically bounded even in the tight-loop
- Data copying on the read-side - its more of an IPC than a lock
- Worst-case: bounded - no lockstep update/read

# pBseq parameters

## ÖSADL

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

Context

Alternatives

**pBseq**

Inherent
randomness

Conclusions

- Number of replicas
- Data array size
- Treatment of index
    - -> collision avoidance (seqcount_latch)
    - -> lockstep avoidance
- Memory model of architecture
    - TSO: no rmb()/wmb() needed (no shared writable data)
    - non-TSO: barriers needed

# Implementation on TSO systems

Writer:

```
for (b = 0; b < NUM_REP; b++) {
  p = &seq.bseq[b];
  p->s++;
  for (d = 0; d < D_SIZE; d++){
    randbit=(~randbit);
    p->data[d] = ...;
    randbit=(~randbit);
  }
  p->s++;
}
```

Reader:

```
static unsigned int idx = 1;
idx -= randbit;
do {
  lp = &seq.b[(--idx)%NUM_REP];
  lseq = lp->s;
  for (d = 0; d < D_SIZE; d++){
    ldata[d] = lp->data[d];
  }
  randbit=(~randbit);
} while (lseq%2 || lseq!=lp->s);
```

# Implementation on non-TSO (ARM/PPC)

- Failure rate is actually quite small (10E-7/call)
- Memory barriers mandatory and expensive
- Works - but no systematic testing yet
- Performance issues and optimization open
- Unclear if a barrier free version is actually possible

Not done yet - the distributions look more or less the same, but there are rare failures if there is a barrier free mitigation is not yet clear - still working on basics here...

# The retry lockstep problem

- Temporary lockstep behavior possible
- Even short lock-stepping could have significant cache side-effects



idle AMD Phenom 2X, int array size=4, temporary lockstep behavior

# Mitigations solution-space

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

Context

Alternatives

pBseq

Inherent
randomness

Conclusions

- Defined writer start conditions and inverted read access direction
- Defined writer start conditions and fixed offsets (seqcount_latch)
- Use of random index (symmetry-breaking)
- Larger replica arrays
- Execution interleaving
- System load

...and any combination of the above

# Defined access direction comparison

OSADL

The distribution indicates that there are very long lockstep sequences possible in tight-loop runs.

# Defined index offset comparison

i7 8-core, cmp (--idx) with (--idx + 1), idle system, data size=4

- unsigned int idx=1; -> lp = &seq.bseq[(−idx)
- unsigned int idx=0; -> lp = &seq.bseq[(−idx + 1)

# Non-determinism

ÖSADL

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

Context

Alternatives

pBseq

Inherent
randomness

Conclusions

- Explicit random numbers
  - Possible but too expensive
- Utilize the asynchronity of system
  - Non-synchronized race on global var
- Utilize the history of the system
  - Static declaration index variable

A reliable entropy source suitable for low-level algorithms to ensure symetry breaking is not only usable for pBseq/pWCS but is a pre-requisite to bring probabilistic solutions to low-level algorithms in general.

# Entropy harvesting - randbits

OSADL

Core of the essed.c

```
unsigned int coin;

Thread A                    Thread B
                              unsigned int draw[2];
  while(active){            drwa[0]=coin;
    coin=~coin;             draw[1]=coin;
                            if(draw[0]<draw[1])
                              inbuf|=(1<<pos);
                            else
                              inbuf&=~(1<<pos);
                            pos++;
```

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
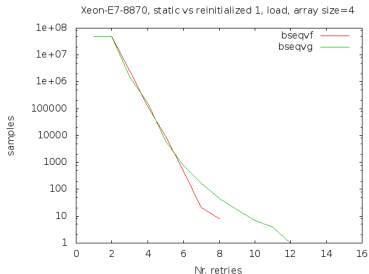Guire
<safety@osadl.

Outline

Context

Alternatives

pBseq

Inherent
randomness

Conclusions

# Entropy harvesting - state history

```
<        static unsigned int idx=1;
---
>        unsigned int idx=1;
```

# Entropy harvesting - state history

ŌSADL

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

Context

Alternatives

pBseq

Inherent
randomness

Conclusions

```
<        static unsigned int idx=1;
---
>        unsigned int idx=1;
```



Xeon-E7-8870, static vs reinitialized 1, load, array size=4

# Entropy harvesting - state history

ӦSADL

pBseq:
probabilistic
buffer-
swapping
sequence
locks

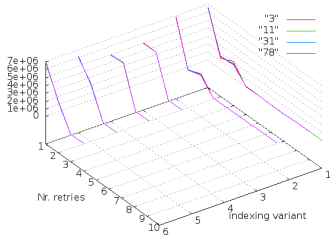Nicholas Mc
Guire
<safety@osadl.

Outline

Context

Alternatives

pBseq

Inherent
randomness

Conclusions

```
<          static unsigned int idx=1;
---
>          unsigned int idx=1;
```

# Entropy insertion points

ÖSADL

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

Context

Alternatives

pBseq

Inherent
randomness

Conclusions

```
int idx  <----- idx = randbit;
         <----- static int idx = 0;
         <----- idx -= randbit;
do {
         lp = &seq.bseq[(--idx)%NUM_REPLICA]; <--- +randbit
         lseq = lp->s;
         for (d = 0; d < DATA_SIZE; d++) {
                 ldata[d] = lp->data[d];
                     <-------------------randbit=(~randbit);
         }
           <--------------------------randbit=(~randbit);
} while (lseq%2 || lseq != lp->s);
```

# Key problem - where to insert entropy

Xeon E7 8870, CPU 0 -> 3,11,31,78, index method cmp, array size 128

# Key problem - where to insert entropy

Looking at the retrydistribution in conjunction with the index distribution

# Worst-case: minimal delays + idle system

work loop impact, load4, i3-3250M int array size=4,32

# Worst-case: minimal delays + idle system

Exactly one worst case scenario -> exhaustive testing possible ?

# Current "best" code

ÖSADL

pBseq:
probabilistic
buffer-
swapping
sequence
locks

Nicholas Mc
Guire
<safety@osadl.

Outline

Context

Alternatives

pBseq

Inherent
randomness

Conclusions

```
static unsigned int idx=1;
for (n = 1; n <= N; ++n) {
        idx += randbit;
        do {
                lp = &seq.bseq[(--idx)%NUM_REPLICA];
                lseq = lp->s;
                for (d = 0; d < DATA_SIZE; d++) {
                        ldata[d] = lp->data[d];
                }
                randbit=(~randbit);
        } while (lseq%2 || lseq != lp->s);
```

# Properties of pBseq

- Freshness: as good as spinlocks/mutexes
- Context separation: reader/writer lock-free/wait-free always
- Performance: Statistically bounded retries
- Worst Case:
  - Idle-system
  - Tight-loop
  - Sall data
  - -> exhaustive testable worst case

# Conclusions

- Related: pWCS and similarities with sequcount_latch
- pBseq is more of an IPC mechanism than a lock (copying semantics)
- Exhaustive testing of worst-case possible
- Simple inherent entropy harvesting is usable in low-level algorithms
- Symmetry breaking can be implemented in a highly reliable manner -> robust guarantees on retries
- Using inherent non-determinism in a systematic form (code wise) is still an open issue
- The code sensitivity to very small changes is large - making the evaluation of the code very hard.

We think its worth digging deeper - synchronization based on robust statistical properties may be a scalable alternative to deterministic locking.