

# Selectively Uniform Concurrency Testing

**Huan Zhao**

Co-Authors: Dylan J. Wolff Umang Mathur Abhik Roychoudhury

{ zhaohuan, wolffd, umathur, abhik } @comp.nus.edu.sg



# Overview

## Background

Model Controlled Concurrency Testing (CCT)  
as a sampling process

## Our Work

What makes a sampling process ideal?

# Concurrent Programs (1)

Critical digital infrastructure



Disastrous concurrency bugs

Northeastern Blackout (2003)



Therac-25 Radiation Therapy Machine (1980s)



Boeing 787 Dreamliner (2013)



# Concurrent Programs (2)

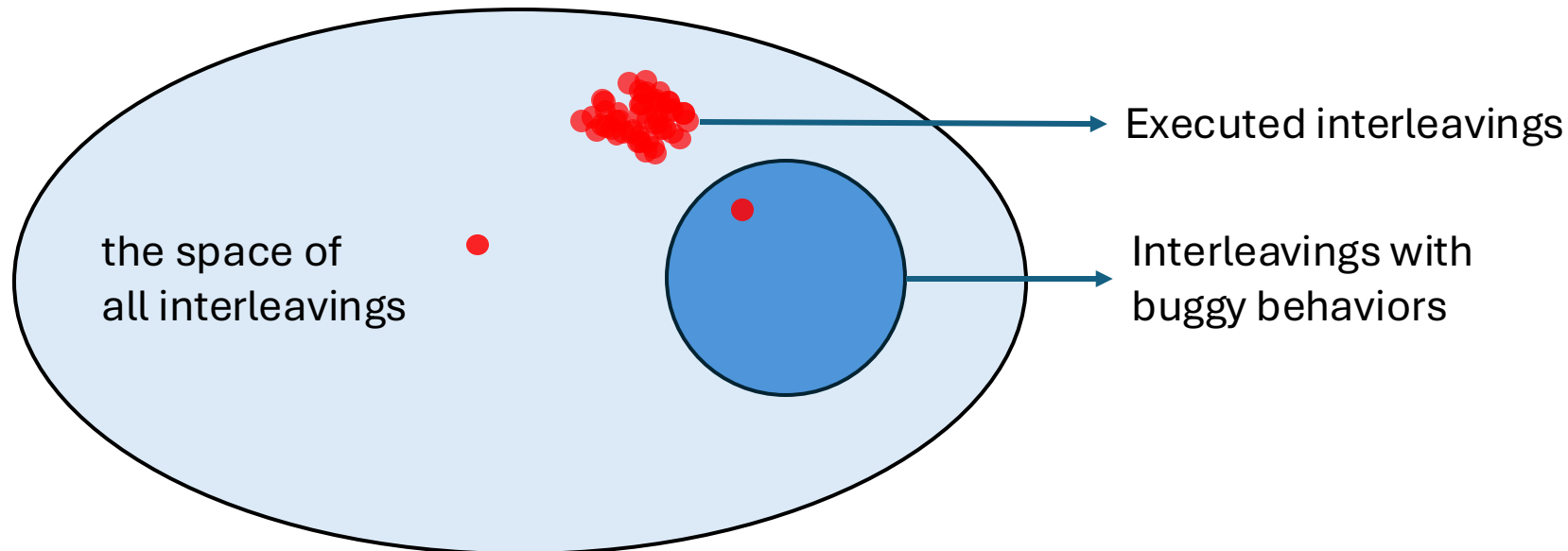
- Non-deterministic behaviors depending on *interleavings*, i.e., ordering of thread executions
  - With 2 threads each with 5 **atomic** events, there are 252 possible observable behaviors (value of x)
- Concurrency bugs may only manifest in few interleavings
- Difficult to expose and reproduce concurrency bugs!

```
1  void thread_A() {  
2      x = x<<1;  
3      x = x<<1;  
4      x = x<<1;  
5      x = x<<1;  
6      x = x<<1;  
7  }
```

```
1  void thread_B() {  
2      x = x<<1+1;  
3      x = x<<1+1;  
4      x = x<<1+1;  
5      x = x<<1+1;  
6      x = x<<1+1;  
7  }
```

# Stress Testing

- Repeated and uncontrolled executions with the default OS scheduler
- However, wildly different interleavings are possible when deployed, due to unpredictable changes to the system workload!

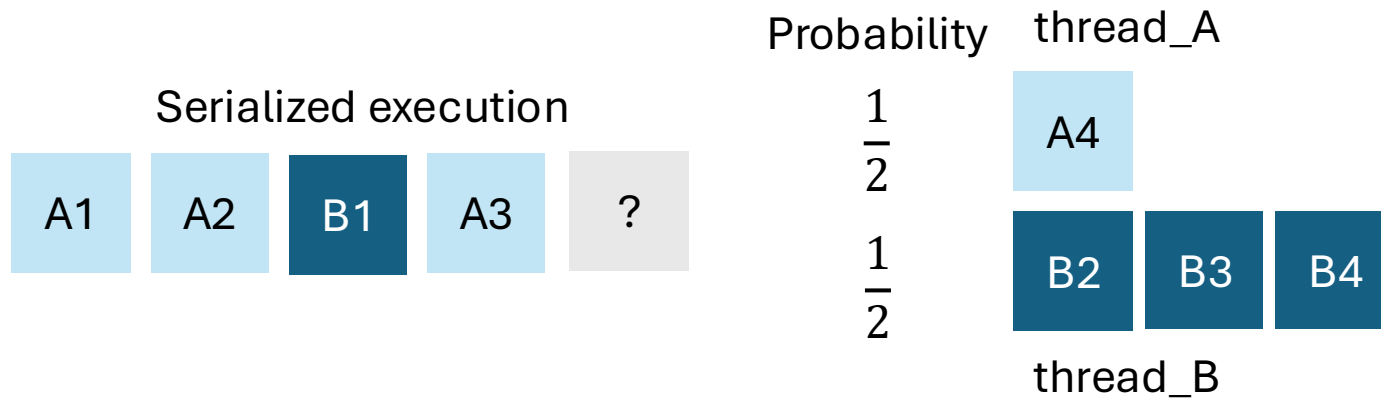


# Controlled Concurrency Testing (CCT)

1. *Serialized* program execution
2. *Online* decision of which event to run next

# Randomized CCT: Random Walk

1. *Serialized* program execution
2. Choose each thread to run *with equal prob.*

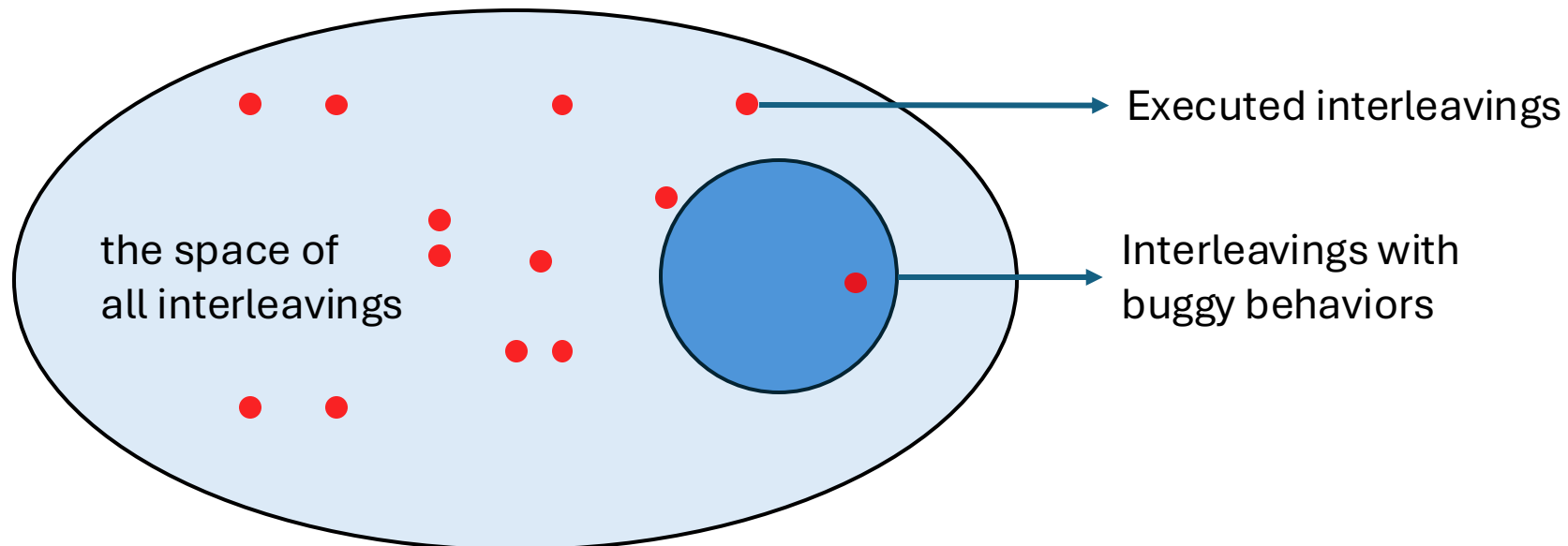


```
1  void thread_A() {  
2      x = x<<1;  
3      x = x<<1;  
4      x = x<<1;  
5      x = x<<1;  
6      x = x<<1;  
7  }
```

```
1  void thread_B() {  
2      x = x<<1+1;  
3      x = x<<1+1;  
4      x = x<<1+1;  
5      x = x<<1+1;  
6      x = x<<1+1;  
7  }
```

# CCT as a Sampling Process

- Sampling of interleavings according to a prob. distribution
- Lightweight yet effective schedule generation, often with prob. guarantees
- Reproduces exposed bugs deterministically



What makes  
a sampling process ideal?

# Uniformity

The **optimal** sampling strategy:

Sample each *program behavior* with equal probability (i.e., uniformly)

→ maximize the *minimum probability* of any behavior being sampled

However, behaviors are program dependent and unknown a priori!

The most fine-grained notion of behaviors: interleavings!

# Motivating Example (1)

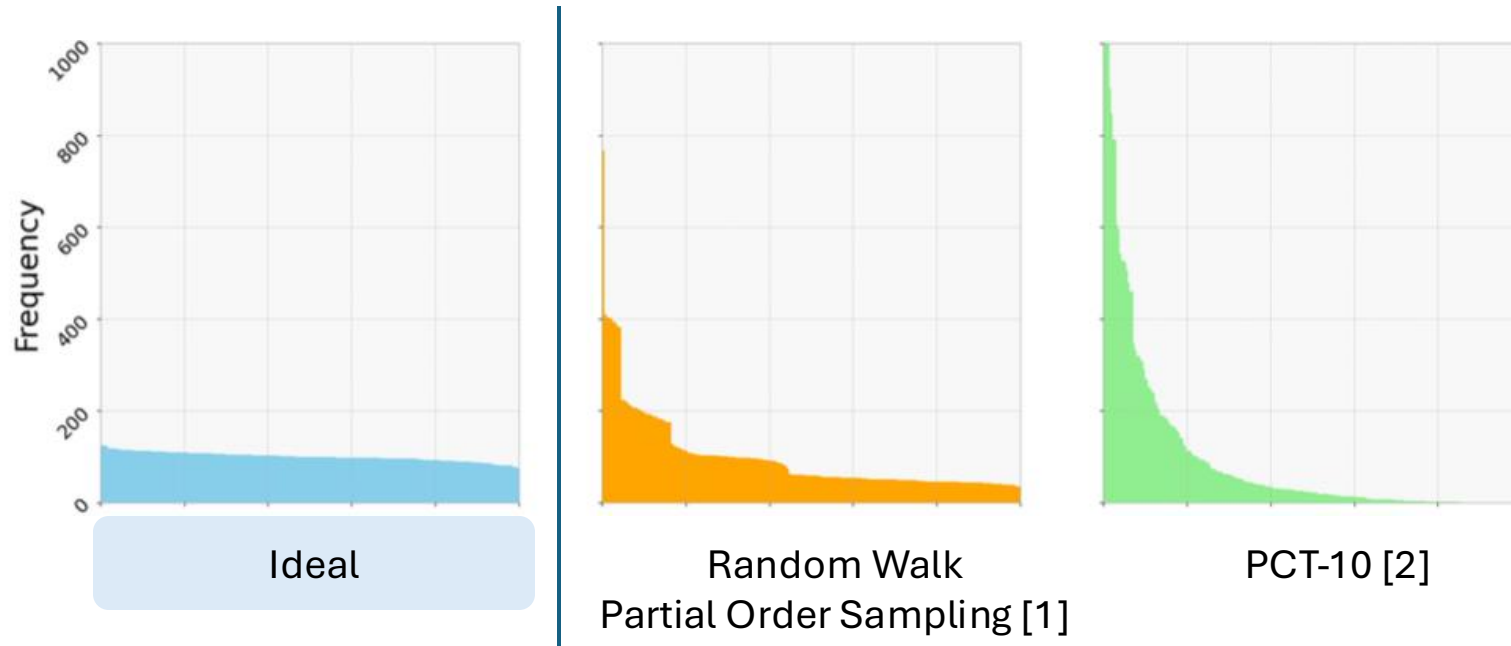
- Interleavings  $\longleftrightarrow$  (observable) program behaviors, i.e., value of  $x$
- Total number of interleavings / behaviors:  $(10 \text{ choose } 5) = 252$
- An ideal algorithm samples each value of  $x$  w.p.  $1/252$

```
1  void thread_A() {  
2      x = x<<1;  
3      x = x<<1;  
4      x = x<<1;  
5      x = x<<1;  
6      x = x<<1;  
7  }
```

```
1  void thread_B() {  
2      x = x<<1+1;  
3      x = x<<1+1;  
4      x = x<<1+1;  
5      x = x<<1+1;  
6      x = x<<1+1;  
7  }
```

# Motivating Example (2)

Total # of int. / behaviors = 252



```
1 void thread_A() {  
2     x = x<<1;  
3     x = x<<1;  
4     x = x<<1;  
5     x = x<<1;  
6     x = x<<1;  
7 }
```

```
1 void thread_B() {  
2     x = x<<1+1;  
3     x = x<<1+1;  
4     x = x<<1+1;  
5     x = x<<1+1;  
6     x = x<<1+1;  
7 }
```

Existing strategies are *highly* biased!

For RW / POS: 2 events are sampled w.p. **1/32**; but 140 events are sampled w.p. **1/512**

For PCT-10: **38** interleavings are not witnessed *once* in >25,000 executions in our experiment!

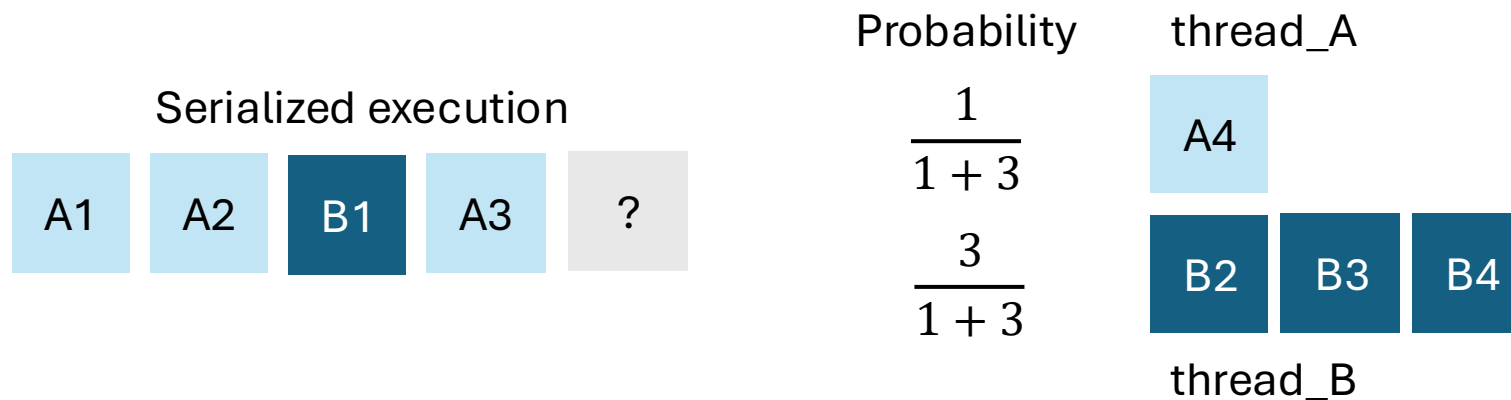
[1] Yuan, Xinhao, et al., 2018. "Partial order aware concurrency sampling." CAV 2018.

[2] Burckhardt, Sebastian, et al., 2010. "A randomized scheduler with probabilistic guarantees of finding bugs." ACM SIGARCH Architecture News 38(1).

## Key Insight 1: Interleaving Uniformity

(a) Interleaving is a proxy of program behaviors;

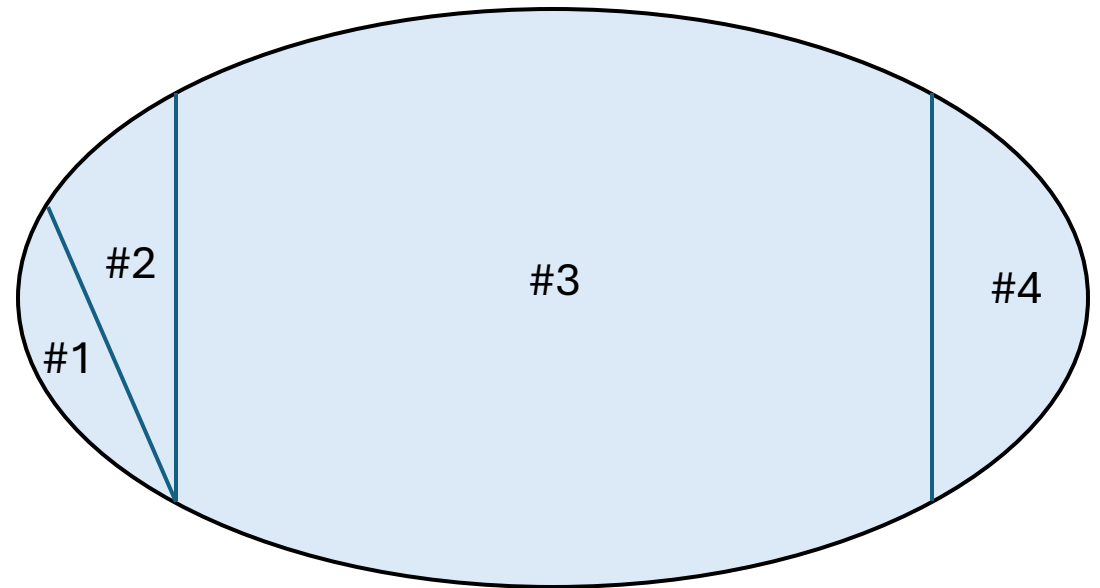
(b) Interleaving uniformity is desirable,  
and achievable with a **weighted** random walk



# Selectivity

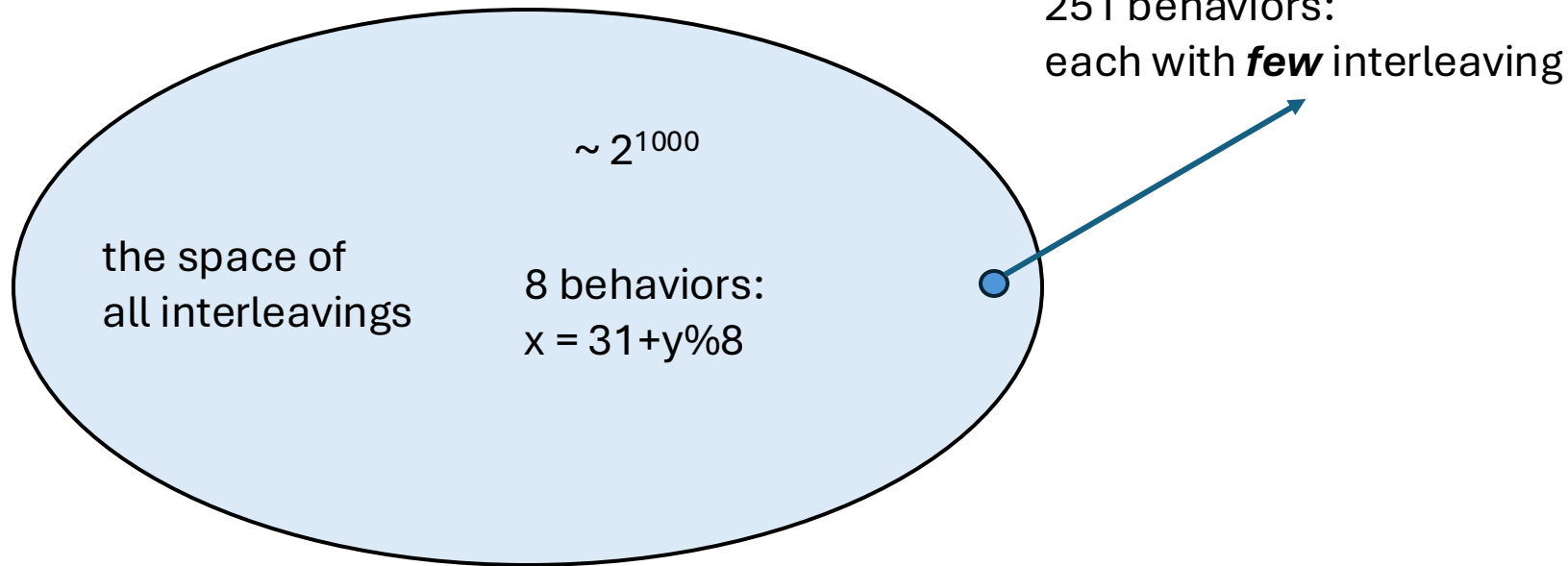
- Interleaving  $\neq$  behavior. Interleaving uniform  $\neq$  behavior uniform.  
e.g., 10 interleavings may result in the same observable behavior (value of x)!
- With interleaving uniformity, most samples result in behavior #3!
- Behaviorally redundant samples!

the space of all interleavings,  
partitioned by *program behaviors*



# Motivating Example (1)

- Interleaving uniformity means that most behaviors will almost never be sampled!

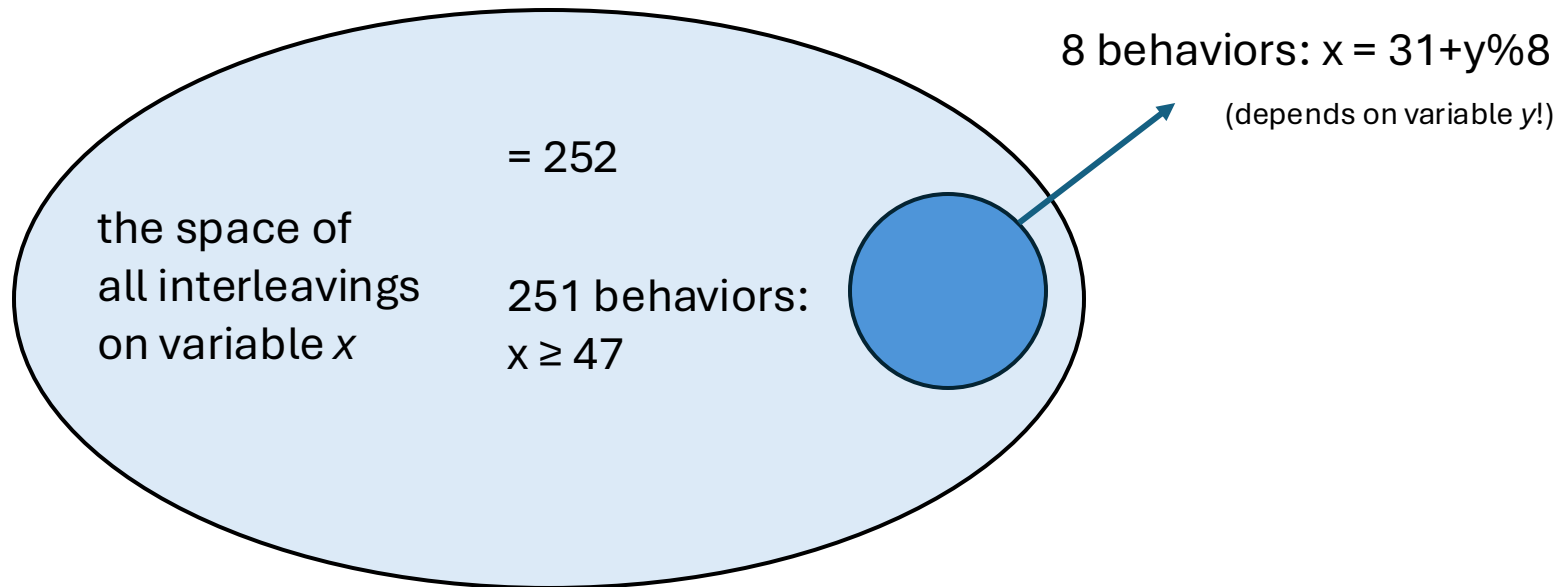


```
1 void thread_A() {  
2     x = x<<1;  
3     x = x<<1;  
4     x = x<<1;  
5     x = x<<1;  
6     x = x<<1;  
7     y = y<<1;  
8     // repeat 1000x  
9     y = y<<1;  
10 }
```

```
1 void thread_B() {  
2     y = y<<1+1;  
3     // repeat 1000x  
4     y = y<<1+1;  
5     x = x<<1+1;  
6     x = x<<1+1;  
7     x = x<<1+1;  
8     x = x<<1+1;  
9     x = x<<1+1+y%8;  
10 }
```

# Motivating Example (2)

- The projected interleavings of variable  $x$  are much more balanced in the behavioral space!



```
1 void thread_A() {  
2     x = x<<1;  
3     x = x<<1;  
4     x = x<<1;  
5     x = x<<1;  
6     x = x<<1;  
7     y = y<<1;  
8     // repeat 1000x  
9     y = y<<1;  
10 }
```

```
1 void thread_B() {  
2     y = y<<1+1;  
3     // repeat 1000x  
4     y = y<<1+1;  
5     x = x<<1+1;  
6     x = x<<1+1;  
7     x = x<<1+1;  
8     x = x<<1+1;  
9     x = x<<1+1+y%8;  
10 }
```

## Key Insight 2: Selective uniformity

Uniform sampling of the interleavings of  
an appropriate subset of program events  
achieves effective *behavioral* exploration

## Key Insight 2: Selective uniformity

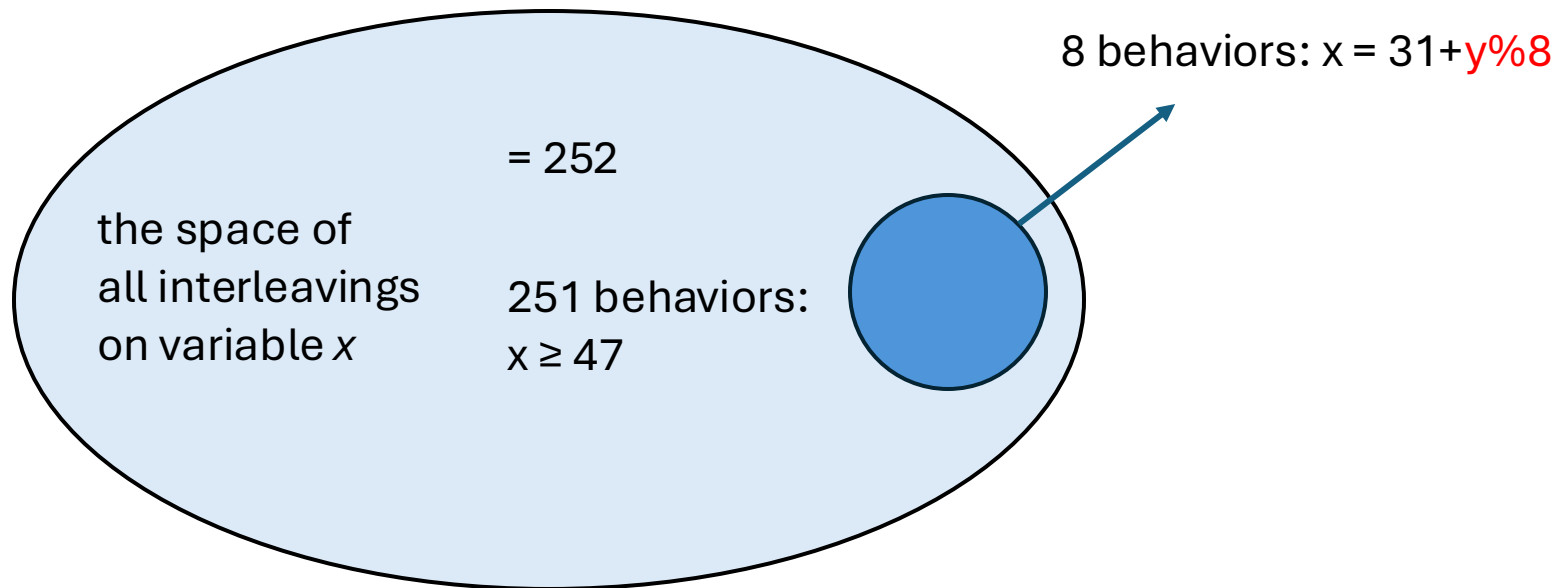
Uniform sampling of the interleavings of  
an appropriate subset of program events  
achieves effective *behavioral* exploration

Given a subset of event  $\Delta \subseteq \Sigma$ , we want to

[Uniformity] sample the interleaving of  $\Delta$  uniformly

# Motivating Example (3)

- At the same time, we should not disable any interleaving of  $y$ , or between  $x$  and  $y$ !



```
1 void thread_A() {
2     x = x<<1;
3     x = x<<1;
4     x = x<<1;
5     x = x<<1;
6     x = x<<1;
7     y = y<<1;
8     // repeat 1000x
9     y = y<<1;
10 }
```

```
1 void thread_B() {
2     y = y<<1+1;
3     // repeat 1000x
4     y = y<<1+1;
5     x = x<<1+1;
6     x = x<<1+1;
7     x = x<<1+1;
8     x = x<<1+1;
9     x = x<<1+1+y%8;
10 }
```

## Key Insight 2: Selective uniformity

Uniform sampling of the interleavings of  
an appropriate subset of program events  
achieves effective *behavioral* exploration

Given a subset of event  $\Delta \subseteq \Sigma$ , we want to

- (1) [Uniformity] sample the interleaving of  $\Delta$  uniformly
- (2) [Completeness] sample any interleaving over  $\Sigma$  with non-zero probability

# Selective Uniform Random Walk

Given a subset of event  $\Delta \subseteq \Sigma$ , we want to

- (1) [Uniformity] sample the interleaving of  $\Delta$  uniformly
- (2) [Completeness] sample any interleaving over  $\Sigma$  with non-zero probability

Central Idea:

- Make *eager* decisions about the interleaving of  $\Delta$  with weighted random walk
- Only schedule  $\Sigma - \Delta$  so that the decision on  $\Delta$  is respected

Constant time per step:  $O(\# \text{ of threads})$

More details discussed in the paper!

---

**Algorithm 2:** SURW

---

```
1 Input:  $\Delta$ ;           // set of interesting events
2 Input:  $n_1, \dots, n_k$ ; // interesting event counts
3  $T_{iNext} \leftarrow$  random  $T_i$  weighted by  $n_i$ ;
4  $blocked \leftarrow \emptyset$ ;
5 while  $E \leftarrow \text{getEnabled}() \neq \emptyset$  do
6    $T_t \leftarrow \text{pickFrom}(E - blocked)$ ;
7   if  $\text{nextEvent}(T_t) \in \Delta$  then
8     if  $T_{iNext} == T_t$  then
9        $n_t \leftarrow n_t - 1$ ;
10       $T_{iNext} \leftarrow$  random  $T_i$  weighted by  $n_i$ ;
11       $blocked \leftarrow \emptyset$ ;
12    else
13       $blocked.add(T_t)$ ; continue;
14   $\text{execute}(\text{nextEvent}(T_t))$ ;
```

---

# Evaluation (1)

[RQ1] Is SURW better at exposing bugs compared to other concurrency testing algorithms?

[RQ2] How do the two key components of SURW, uniformity and selectivity, contribute to its effectiveness?

# Evaluation (2)

3 established concurrency testing benchmarks in the community [1-3]

Other state-of-the-art algorithms: PCT-3 [4], PCT-10 [4], POS [5]

Baselines: Random Walk, Non-Selective, Non-Uniform

[1] Thomson, Paul, et al., 2016. "Concurrency testing using controlled schedulers: An empirical study." *TOPC* 2.4 (2016).

[2] Meng, Ruijie, et al., 2019. "ConVul: an effective tool for detecting concurrency vulnerabilities." *ASE* 2019.

[3] Liang, Jiashuo, et al., 2023. "RaceBench: A Triggerable and Observable Concurrency Bug Benchmark." *AsiaCCS* 2023.

[4] Yuan, Xinhao, et al., 2018. "Partial order aware concurrency sampling." *CAV* 2018.

[5] Burckhardt, Sebastian, et al., 2010. "A randomized scheduler with probabilistic guarantees of finding bugs." *ACM SIGARCH Architecture News* 38(1).

# RQ1 Bug Finding (1)

[Metric 1]. Average # of bugs exposed by different algorithms (the higher, the better)

Benchmark	SURW	PCT	POS	Random Walk
SCTBench (max. 37)	34.90	30.75	29.25	19.90
		+13%	+19%	+75%
RaceBenchData (max. 1,500)	944	461	885	489
		+105%	+7%	+93%

[Metric 2]. Average # of interleavings sampled to expose each bug (the lower, the better)

- On 26 / 35 targets, SURW requires the minimum number of schedules
- On 6 / 9 other targets, SURW requires only < 10 schedules on average

# RQ1 Bug Finding (2)

Target	SURW	PCT-3	PCT-10	POS	Random Walk
CS/twostage	$8 \pm 4$	$13 \pm 14$	$9 \pm 10$	$15 \pm 12$	$464 \pm 581$
CS/twostage_20	$6 \pm 3$	$159 \pm 151$	$101 \pm 92$	$156 \pm 137$	—
CS/twostage_50	$20 \pm 17$	$1676 \pm 1715$	$692 \pm 392$	$1637 \pm 1385$	—
CS/twostage_100	$454 \pm 444$	$7466 \pm 831^*$	$5726 \pm 2591^*$	$6674 \pm 2877^*$	—
CS/reorder_3	$7 \pm 7$	$185 \pm 199$	$148 \pm 191$	$86 \pm 70$	—
CS/reorder_4	$7 \pm 6$	$554 \pm 643$	$362 \pm 213$	$533 \pm 651$	—
CS/reorder_5	$10 \pm 9$	$647 \pm 517$	$1094 \pm 1216$	$2169 \pm 2182$	—
CS/reorder_10	$17 \pm 11$	$3225 \pm 2426^*$	$4462 \pm 3266^*$	—	—
CS/reorder_20	$6 \pm 4$	$3005 \pm 2680$	$3297 \pm 2877^*$	—	—
CS/reorder_50	$13 \pm 12$	$3304 \pm 1721^*$	—	—	—
CS/reorder_100	$194 \pm 214$	—	—	—	—
ConVul/CVE-2013-1792	$15 \pm 13$	$95 \pm 83$	$50 \pm 61$	$39 \pm 33$	$364 \pm 289$
ConVul/CVE-2016-1972	$11 \pm 8$	$4902 \pm 2391^*$	$2712 \pm 2704^*$	$34 \pm 34$	$299 \pm 256$
ConVul/CVE-2016-1973	$5 \pm 3$	$10 \pm 8$	$6 \pm 3$	$5 \pm 4$	$308 \pm 333$
ConVul/CVE-2016-7911	$8 \pm 9$	$20 \pm 18$	$15 \pm 15$	$11 \pm 9$	$3 \pm 2$
ConVul/CVE-2016-9806	$3 \pm 2$	$7 \pm 6$	$4 \pm 3$	$7 \pm 5$	$2209 \pm 2065$
ConVul/CVE-2017-15265	—	—	—	—	—
ConVul/CVE-2017-6346	$15 \pm 10$	$24 \pm 18$	$20 \pm 19$	$10 \pm 9$	$3 \pm 4$

CS/twostage\_100:  
SURW in  $\sim 450$  schedules  
vs.  $>5k$  schedules

CS/reorder\_100:  
SURW in  $\sim 200$  schedules  
vs.  $>> 200k$  schedules

CVE-2016-1972:  
SURW in  $\sim 10$  schedules  
vs. PCT  $>2k$  schedules

# RQ1 Bug Finding (2)

## Answer to RQ1

SURW outperforms other sampling algorithms by a large margin in bug finding!

# RQ2 Ablation Study (1)

Two ablative versions:

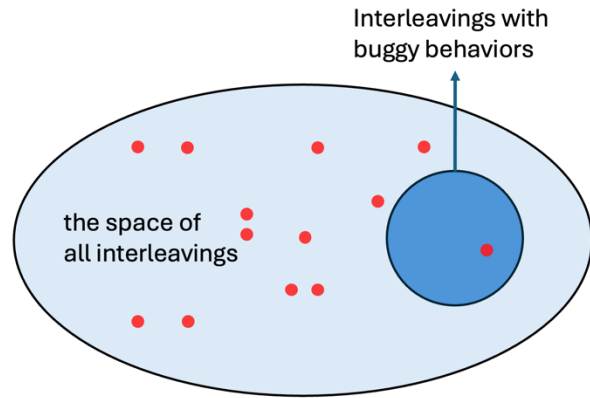
- Non-Uniform (N-U): Selective interesting subset + naïve random walk
- Non-Selective (N-S): Weighted random walk on the entire program

[Metric 1]. Average # of bugs exposed by different algorithms

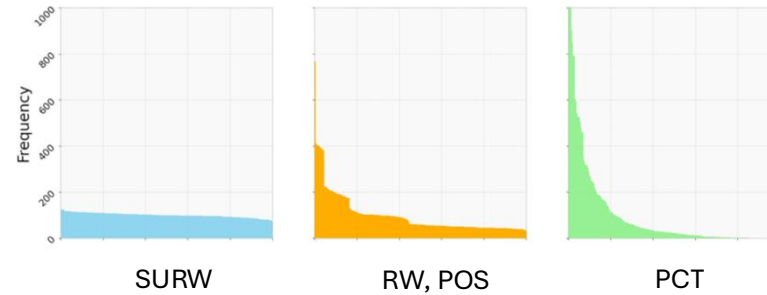
Benchmark	<b>SURW</b>	Non-Uniform	Non-Selective
SCTBench (max. 37)	<b>34.90</b>	29.70	30.75
		+18%	+14%

# Summary

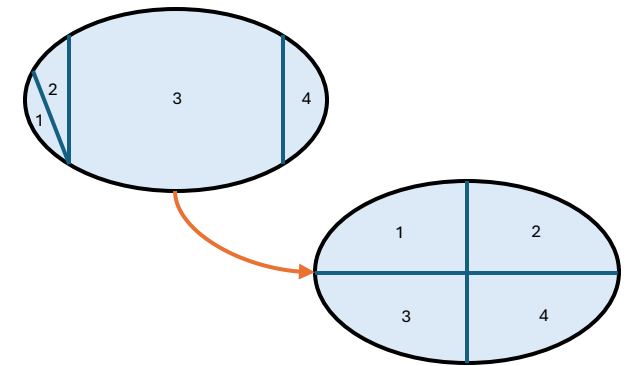
## Concurrency testing as sampling



## Interleaving-uniformity sampling via weighted random walk



## Effective behavioral exploration via selective uniformity



Link to Artifact

Concurrency testing as sampling

**selectively uniform**



Link to Paper