

A Support for Persistent Memory in Java

Anatole Lefort

Advisors: Pierre Sutra, Gaël Thomas

Télécom SudParis Institut Polytechnique de Paris

Ph.D. Defense – Mar. 24th, 2023



A Support for Persistent Memory in Java

Anatole Lefort

Advisors: Pierre Sutra, Gaël Thomas

Télécom SudParis Institut Polytechnique de Paris

Ph.D. Defense – Mar. 24th, 2023

Data Persistence

Background

Computer programs deal with two kinds of data.

- Transient:

- Limited Lifetime: renewed at every program execution.
 - do not survive crashes.
- hosted in Main Memory

- Persistent:

- Extended Lifetime: recalled and reused in subsequent executions.
 - remain consistent even in the wake of a failure.
- hosted on Storage Devices

"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms



"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms



"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms



"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms



"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms



"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms



"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms



"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

Common media: HDD, Flash disks

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms

1 - Dual data representation
 2 - Expensive I/Os



"Persistence is the ability of data to outlive the instance of a program"

"Data remain persistent for the extent of time during which they may be recalled and used by a program"

Common media: HDD, Flash disks

- **Durable**: resist reboots, power loss
- Large capacity: at least TBs
- Slow access latency: 200µs-15ms

1 - Dual data representation

2 - Expensive I/Os



Background

PMEM: A memory device on which data survives power cycles.

What changes with PMEM ?



- 1 Dual data representation
- 2 Expensive I/Os

Background

PMEM: A memory device on which data survives power cycles.

What changes with PMEM ? 1- add PMEM DIMM I - add PMEM DIMM (Persistent) What changes with PMEM ? PMEM PMEM PMEM (Volatile) (Volatile)

- 1 Dual data representation
- 2 Expensive I/Os

Background

PMEM: A memory device on which data survives power cycles.

What changes with PMEM ? 1- add PMEM DIMM, CPU-attached persistent media PMEM PMEM PMEM CPU DRAM (Volatile)

- 1 Dual data representation
- 2 Expensive I/Os

Background

PMEM: A memory device on which data survives power cycles.



1 - Dual data representation

2 - Expensive I/Os

Background

PMEM: A memory device on which data survives power cycles.

What changes with PMEM ?

 add PMEM DIMM, CPU-attached persistent media
 move persistent data, disks become redundant



1 - Dual data representation

2 - Expensive I/Os

Background

PMEM: A memory device on which data survives power cycles.

What changes with PMEM ?

 add PMEM DIMM, CPU-attached persistent media
 move persistent data, disks become redundant
 no more disk I/Os



1 - Dual data representation 2 - Expensive I/Os

Background

PMEM: A memory device on which data survives power cycles.

What changes with PMEM ?

 add PMEM DIMM, CPU-attached persistent media
 move persistent data, disks become redundant
 no more disk I/Os
 working copy of data is durable



1 - Dual data representation 2 - Expensive I/Os

Background

PMEM: A memory device on which data survives power cycles.

What changes with PMEM ?

 add PMEM DIMM, CPU-attached persistent media
 move persistent data, disks become redundant
 no more disk I/Os
 working copy of data is durable





Background

PMEM: A memory device on which data survives power cycles.

What changes with PMEM ? 1- add PMEM DIMM, CPU-attached persistent media 2- move persistent data, disks become redundant 3- no more disk I/Os 4- working copy of data is durable



Benefits

- 1- No more (un)marshalling
- 2- No need for data caching
- 3- Faster recovery
- 4- Lower software complexity



new persistent medium (in-between SSD and DRAM)

Durable

resists reboots, power loss

High-density smallest DIMM = 128 GB up to 8x DDR4 capacity

Byte addressable persistent memory abstraction

High-performance

low latency (seq. read/write ~ 160/90ns) high bandwidth (up to 8.10GB/s, *2nd gen*)



Intel Optane PMEM, 2019

Background

Non-volatile main memory (NVMM)

Background



Non-volatile main memory (NVMM)





Non-volatile main memory (NVMM)

Background



Background

How do we use it ? Storage device compatibility mode (1) ? Persistent Memory (2) ?

(1) File system interface

open/close/read/write/sync

(2) Direct memory access mmap



Intel Optane PMEM, 2019

Background

How do we use it ? Storage device compatibility mode (1) ?

(1) File system interface

open/close/read/write/sync



Background

How do we use it ? Storage device compatibility mode (1) ?

(1) File system interface

open/close/read/write/sync



Disabling durability significantly boost performance
Dummy file systems are seemingly identical to a PMEM FS

Background

How do we use it ? Storage device compatibility mode (1) ?

(1) File system interface

open/close/read/write/sync



Software Bottlenecks:

- dual representation (consistency)
- cost of marshalling

Disabling durability significantly boost performance
Dummy file systems are seemingly identical to a PMEM FS

Background

How do we use it ? Persistent Memory (2) ?

(2) Direct memory access mmap

+ memory load/store CPU instructions

Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access mmap
 - + memory load/store CPU instructions



Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access mmap
 - + memory load/store CPU instructions



Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access mmap
 - + memory load/store CPU instructions

data update path 1- pull data in CPU caches



Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access mmap
 - + memory load/store CPU instructions

<u>data update path</u> 1- pull data in CPU caches 2- update data in CPU caches first,



Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access mmap
 - + memory load/store CPU instructions

<u>data update path</u> 1- pull data in CPU caches 2- update data in CPU caches first, NVMM version becomes stale



Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access mmap
 - + memory load/store CPU instructions



data update path

1- pull data in CPU caches 2- update data in CPU caches first, NVMM version becomes stall 3- Hardware dictates when and in which

3- Hardware dictates when and in which order data are evicted from caches
Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access mmap
 - + memory load/store CPU instructions



data update path

 1- pull data in CPU caches
 2- update data in CPU caches first, NVMM version becomes stall
 3- Hardware dictates when and in which order data are evicted from caches
 4- update eventually reaches NVMM

Background

How do we use it ? Persistent Memory (2) ?

(2) Direct memory access mmap

+ memory load/store CPU instructions



data update path

 1- pull data in CPU caches
 2- update data in CPU caches first, NVMM version becomes stall
 3- Hardware dictates when and in which order data are evicted from caches
 4- update eventually reaches NVMM

Version in caches (<>) diverged from NVMM (<>)

Is it always **safe** to recover (**◇**) **after a crash** ? <u>How to keep **data crash-consistent** ?</u>

Background

How do we use it ? Persistent Memory (2) ?

(2) Direct memory access mmap

+ memory load/store CPU instructions



+ special **flush/fence** CPU instructions (manually control cache line eviction order)



Background

How do we use it ? Persistent Memory (2) ?

(2) **Direct memory access** mmap

+ memory **load/store** CPU instructions



data update path

1- pull data in CPU caches 2- update data in CPU caches first, NVMM version becomes stall 3- Hardware dictates when and in which order data are evicted from caches 4- update eventually reaches NVMM

- + special **flush/fence** CPU instructions (manually control cache line eviction order)
- Too low-level programming \Rightarrow
- Brittle reasoning about crash-consistency \Rightarrow



Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access ~ <u>the easy way</u> mmap
 - + fitting programming abstractions (e.g. Intel's PMDK)
 - ⇒ ensure data **crash-consistency**
 - ⇒ aid data **recovery**

Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access ~ <u>the easy way</u> mmap
 - + fitting programming abstractions (e.g. Intel's PMDK)
 - ⇒ ensure data **crash-consistency**
 - ⇒ aid data **recovery**

<u>challenges:</u>

- _ support for **failure-atomicity** abstractions ?
- _ persistent memory allocator ?
- _ persistent pointers in ephemeral process address space ?

Background

How do we use it ? Persistent Memory (2) ?

- (2) Direct memory access ~ <u>the easy way</u> mmap
 - + fitting programming abstractions (e.g. Intel's PMDK)
 - ⇒ ensure data **crash-consistency**
 - ⇒ aid data **recovery**

<u>challenges:</u>

- _ support for **failure-atomicity** abstractions ?
- _ persistent memory allocator ?
- _ persistent pointers in ephemeral process address space ?

Language-level NVMM programming help & support

Managed language - first released in the 90s - still an industry standard.

Many data stores & processing frameworks:

- Spark, Hadoop, Kafka, Flink, Cassandra, HBase, Elasticsearch, etc.



Nearly all NVMM libraries and tools support only <u>native code</u> (C, C++)

Lack of *efficient* interfaces :

- (1) File System [ext4-dax]
 - storage device compatibility mode (cf. slide 6)

(2) Intel's PMDK through the Java Native Interface (JNI) [PCJ]

- native library with compatibility layer
 - slower than **FS** on YCSB benchmark (*cf. evaluation*)

Problem statement: No Java-native NVMM interface

= [Espresso, AutoPersist, go-pmem]

Managed persistent Java objects

= extend JVM to manage persistent memory

Motivation

= [Espresso, AutoPersist, go-pmem]

Managed persistent Java objects

= extend JVM to manage persistent memory



Motivation





Motivation



Contribution Overview

J-NVM

J-NVM - Off-Heap Persistent Objects

Challenges

(1) single data representation

Memory

(2) direct access to NVMM

(3) crash-consistency

Java

(4) object-oriented idioms

(5) explicit persistent types

(6) persistent memory management with low overhead and scalable to large heaps

Contribution Overview

J-NVM - Off-Heap Persistent Objects

Persistent Memory

Java

(1) single data representation

Challenges

(2) direct access to NVMM

(3) crash-consistency

Features

(1, 6) off-heap persistent objects

(2) sun.misc.Unsafe

(3) failure-atomic blocks + fine-grained

(4) object-oriented idioms

(5) explicit persistent types

(6) persistent memory management with low overhead and scalable to large heaps (4) persistent java objects + PDT library

(5) **class-centric** model + code generator

(6) recovery-time GC (no online GC) explicit free()

A Simple Bank:

Server	Bank	Account
	-accounts: Map <string, account=""></string,>	<i>-id</i> : Integer <i>-balance</i> : Long
	+performTransfer(String <i>from</i> , String <i>to</i> , long <i>amount</i>) +createAccount(String <i>id</i> , long <i>initialDeposit</i>)	+transferTo(Account dest, long amount)

Demo - *It's showtime* !

J-NVM

natole@latitude ~/Documents/phd/invm-demo \$ git checkout invm-variant	Transferring \$13966 from 20790 to 25979 OK
Switched to branch 'invm-variant'	Transferring \$807 from 19797 to 17432 0K
Your branch is up to date with 'origin/jnvm-variant'.	Transferring \$26127 from 13282 to 14515 OK
<pre>nnatole@latitude ~/Documents/phd/jnvm-demo \$ mvn clean install -Dmaven.test.sk</pre>	Transferring \$20891 from 15389 to 16612 OK
.p=true	Transferring \$19731 from 25022 to 30933 OK
INFO] Scanning for projects	Transferring \$465 from 16948 to 163 OK
INFO]	Transferring \$14739 from 27212 to 31897 OK
<pre>INFO] eu.telecomsudparis.jnvm:jnvm-demo ></pre>	Transferring \$21187 from 19167 to 6331 OK
	Transferring \$29329 from 2542 to 5080 OK
INFO] Building jnvm-demo 1.0-SNAPSHOT	Transferring \$22303 from 7180 to 7857 OK
INF0][jar]	Transferring \$11984 from 3348 to 31671 OK
	Transferring \$31963 from 11914 to 5062 OK
INF0]	Transferring \$2761 from 16502 to 10200 OK
[INFO] maven-clean-plugin:2.5:clean (default-clean) @ jnvm-demo	Transferring \$8826 from 14802 to 5272 OK
INFOJ Deleting /home/anatole/Documents/phd/jnvm-demo/target	Transferring \$16226 from 11690 to 12212 OK
INFO]	Transferring \$13410 from 24/74 to 27075 UK
[INFU] maven-resources-plugin:2.6:resources (detault-resources) @ jnvm-dem	Transferring \$18111 from 19755 to 3585 UK
TNEO1 Using UNTE 01 anoshing to conv filtered recourses	Transferring \$31013 from 13963 to 26681 UK
INFO] USing UP-6 encouring to copy fittered resources.	Transferring \$22005 from 21501 to 12020 UK
	Transferring \$28280 from 20578 to 12031 OK
TNFO] mayen-compiler-plugin:3 6 l:compile (default-compile) @ ipym-demo -	Transferring \$5633 from 9057 to 21579 OK
rinoj - moven compret pedgritotorreompret (derdate compret) e jivin delp	Transferring \$15372 from 18749 to 27620 OK
	esferring \$30340 from 29898 to 25940 OK
	rsferring \$18655 from 11866 to 3223 OK
	Transferring \$1096 from 22652 to 29958 OK
	Transferring \$20332 from 19758 to 10406 OK
	Transferring \$16902 from 14992 to 26568 OK
	Transferring \$23650 from 17869 to 25875 OK
	Transferring \$28326 from 26926 to 4780 OK
	Transferring \$18147 from 20449 to 10147 OK
	Transferring \$8875 from 20751 to 117 OK
	Transferring \$13754 from 28717 to 30340 OK
	Transferring \$29041 from 18920 to 26579 OK
	Transferring \$12721 from 10616 to 12903 OK
	Transferring \$7333 from 17024 to 5286 0K
	Transferring \$7783 from 1402 to 18889 UK
	Transferring \$11376 from 30535 to 19655 UK
	Transferring \$25517 from 13929 to 4160 UK
	Transferring \$20489 from 6400 to 10204 OK
	Transferring \$21257 from 2044 to 12741 OK
	Transferring \$28362 from 22548 to 30334
	anatole@latitude ~/Documents/hhd/invm-demo \$ /hip/client sh tot
	anatole@latitude ~/Documents/phd/invm-demo \$
4] 0:java*	"latitude" 1

1:04 03-Jun-2

Outline

(1) Introduction

- data persistence
- persistent memory
- NVMM
- why Java?
- prior art
- contribution overview

- (2) System Design of J-NVM
 - demo
 - key idea
 - programming model
 - persistent objects
 - code generator
 - J-PFA
 - J-PDT
- (3) Evaluation
 - YCSB benchmark
 - recovery

(4) Conclusion



J-NVM = Off-Heap Persistent Objects

Key idea each persistent object is *decoupled* into

- <u>a persistent data structure</u>: **unmanaged**, allocated off-heap (NVMM)
- <u>a proxy</u>: **managed**, allocated on-heap (DRAM)

Persistent object is

- <u>a persistent data structure</u>
 - holds object fields
- <u>a proxy</u>
 - holds object methods
 - implements PObject interface
 - intermediates access to pers. data structure
 - instantiated lazily (low GC pressure)

Alive when reachable (from persistent root)

Class-centric model

- safe references thanks to the type system



Sys Design

```
Map root = JNVM.root();
Bank b = root.get("Bank");
Account a = b.find("toto");
a.setBalance(42);
```



- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed

Account a = new Account("toto", 42);

(DRAM)

(NVMM)



- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed

Account a = new Account("toto", 42);

(DRAM)

(NVMM)



- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed

Account a = new Account("toto", 42);

(DRAM)



(NVMM)





- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed

Account a = new Account("toto", 42);





- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed

Account a = bank.find("toto");





- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed

Account a = bank.find("toto");





- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed

Account a = bank.find("toto");





- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed





- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed





(DRAM)



- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed







(NVMM)



- allocate NVMM
- attach persistent data structure

Re-Constructor

- re-attach proxy
- re-build soft state via resurrect()

Destructor

- explicit JNVM.free() to reclaim NVMM
- detach proxy
- ready to be GCed

JNVM.free(a);

(DRAM)

(NVMM)

Overview

J-NVM = Off-Heap Persistent Objects

Tooling

- built-in off-heap memory management for NVMM
- code generator: automatic decoupling for POJOs
- J-PFA: automatic failure-atomic code
- J-PDT: data types + collections for persistent memory
- low-level API (for experts)
- recovery-time GC

Programming model - *code generator*

Implementation

Goals

- compute class-wide off-heap layout
- (2) (3) replace (non-transient) field accesses
- generate constructor, re-constructor
- FA-wrap non-private methods

```
@Persistent(fa="non-private")
class Account {
  PString name;
  int balance;
  transient int y;
  Account(String name, int balance) {
    this.name = new PString(id);
    this.balance = balance;
  }
  void transferTo(Account dest, int amount) {
    this.balance -= amount;
    dest.balance += amount;
```

Programming model - *code generator*

Implementation

Goals

- compute class-wide off-heap layout
- (2) (3) replace (non-transient) field accesses
- generate constructor, re-constructor
- FA-wrap non-private methods

```
@Persistent(fa="non-private")
class Account {
  PString name;
  int balance;
  transient int y;
  Account(String name, int balance) {
    this.name = new PString(id);
    this.balance = balance;
  }
  void transferTo(Account dest, int amount) {
    this.balance -= amount;
    dest.balance += amount;
```

Programming model - code generator

Implementation

Goals



compute class-wide off-heap layout replace (non-transient) field accesses generate constructor, re-constructor FA-wrap non-private methods



```
@Persistent(fa="non-private")
class Account {
    PString name;
    int balance;
    transient int y;
```

```
Account(String name, int balance) {
  this.name = new PString(id);
  this.balance = balance;
```

```
void transferTo(Account dest, int amount) {
  this.balance -= amount;
  dest.balance += amount;
}
```
Implementation

Goals

(1)

(2) (3)

4

compute class-wide **off-heap layout** (a) remove persistent attributes replace (non-transient) field accesses generate constructor, re-constructor FA-wrap non-private methods

(DRAM)



```
// transformed code (decompiled)
class Account {
  transient int y;
 Account(String name, int balance) {
    this.name = new PString(id);
    this.balance = balance;
 void transferTo(Account dest, int amount) {
   this.balance -= amount;
    dest.balance += amount;
```

Implementation

Goals

compute class-wide off-heap layout

 (a) remove persistent attributes
 (b) add "addr" field

 replace (non-transient) field accesses
 generate constructor, re-constructor
 FA-wrap non-private methods

(DRAM)

}



```
// transformed code (decompiled)
class Account {
   long addr; // persistent data structure
   transient int y;
```

```
Account(String name, int balance) {
  this.name = new PString(id);
  this.balance = balance;
```

```
void transferTo(Account dest, int amount) {
   this.balance -= amount;
   dest.balance += amount;
}
```

Implementation

Goals

compute class-wide off-heap layout

 (a) remove persistent attributes
 (b) add "addr" field
 (c) generate or transform field getters/setters

 (2) replace (non-transient) field accesses
 (3) generate constructor, re-constructor
 (4) FA-wrap non-private methods

```
// transformed code (continued)
PString getName() {
   return JNVM.readPObject(addr, 0);
}
protected void setName(PString v) {
   JNVM.writePObject(addr, 0, v);
}
int getBalance() {
   return JNVM.readInt(addr, 8);
}
```

```
}
void setBalance(int v) {
   JNVM.writeInt(addr, 8, v);
```

```
// transformed code (decompiled)
class Account {
   long addr; // persistent data structure
   transient int y;
```

```
Account(String name, int balance) {
  this.name = new PString(id);
  this.balance = balance;
```

```
void transferTo(Account dest, int amount) {
  this.balance -= amount;
  dest.balance += amount;
}
```

Implementation

Goals

. . .

compute class-wide off-heap layout

 (a) remove persistent attributes
 (b) add "addr" field
 (c) generate or transform field getters/setters

 (2) replace (non-transient) field accesses
 (3) generate constructor, re-constructor
 (4) FA-wrap non-private methods

```
// transformed code (continued)
PString getName() {
   return JNVM.readPObject(addr, 0);
}
protected void setName(PString v) {
   JNVM.writePObject(addr, 0, v);
}
int getBalance() {
   return JNVM.readInt(addr, 8);
}
void setBalance(int v) {
   JNVM.writeInt(addr, 8, v);
}
```

```
// transformed code (decompiled)
class Account implements PObject {
   long addr; // persistent data structure
   transient int y;
```

```
Account(String name, int balance) {
  this.name = new PString(id);
  this.balance = balance;
```

```
void transferTo(Account dest, int amount) {
  this.balance -= amount;
  dest.balance += amount;
}
```

1- Use JNVM static helpers

Implementation

Goals

. . .

compute class-wide off-heap layout

 (a) remove persistent attributes
 (b) add "addr" field
 (c) generate or transform field getters/setters

 (2) replace (non-transient) field accesses
 (3) generate constructor, re-constructor
 (4) FA-wrap non-private methods

```
// transformed code (continued)
  PString getName() {
    return JNVM.readPObject(addr, 0);
  protected void setName(PString v) {
    JNVM.writePObject(addr, 0, v);
  int getBalance() {
    return JNVM.readInt(addr, 8);
  void setBalance(int v) {
    JNVM.writeInt(addr, 8, v);
```

```
// transformed code (decompiled)
class Account implements PObject {
   long addr; // persistent data structure
   transient int y;
```

```
Account(String name, int balance) {
  this.name = new PString(id);
  this.balance = balance;
```

```
void transferTo(Account dest, int amount) {
  this.balance -= amount;
  dest.balance += amount;
}
```

1- Use JNVM static helpers with field offsets

Implementation

Goals

. . .

compute class-wide off-heap layout

 (a) remove persistent attributes
 (b) add "addr" field
 (c) generate or transform field getters/setters

 (2) replace (non-transient) field accesses
 (3) generate constructor, re-constructor
 (4) FA-wrap non-private methods

```
// transformed code (continued)
    PString getName() {
        return JNVM.readPObject(addr, 0);
    }
    protected void setName(PString v) {
        JNVM.writePObject(addr, 0, v);
    }
    int getBalance() {
        return JNVM.readInt(addr, 8);
    }
    void setBalance(int v) {
```

```
JNVM.writeInt(addr, 8, v);
```

```
// transformed code (decompiled)
class Account implements PObject {
   long addr; // persistent data structure
   transient int y;
```

```
Account(String name, int balance) {
  this.name = new PString(id);
  this.balance = balance;
```

```
void transferTo(Account dest, int amount) {
  this.balance -= amount;
  dest.balance += amount;
}
```

Use JNVM static helpers with field offsets
 internal setter for *final* fields

Implementation

Goals

compute class-wide off-heap layout

 (a) remove persistent attributes
 (b) add "addr" field
 (c) generate or transform field getters/setters

 (2) replace (non-transient) field accesses
 (3) generate constructor, re-constructor
 (4) FA-wrap non-private methods

```
// transformed code (decompiled)
class Account implements PObject {
   long addr; // persistent data structure
   transient int y;
   Account(String name, int balance) {
    this.setName(new PString(id));
    this.setBalance(balance);
   }
   void transferTo(Account dest, int amount) {
    this.setBalance(getBalance() - amount);
   }
}
```

dest.setBalance(dest.getBalance() + amount);

. . .

Implementation

Goals

```
    compute class-wide off-heap layout

            (a) remove persistent attributes
            (b) add "addr" field
            (c) generate or transform field getters/setters

    (2) replace (non-transient) field accesses
    (3) generate constructor, re-constructor
    (4) FA-wrap non-private methods
```

```
// transformed code (continued)
Account(long addr) {
   this.addr = addr;
   this.resurrect();
}
```

```
// transformed code (decompiled)
class Account implements PObject {
   long addr; // persistent data structure
   transient int y;
```

```
Account(String name, int balance) {
   this.addr = JNVM.alloc(getClass(), size());
   this.setName(new PString(id));
   this.setBalance(balance);
}
```

void transferTo(Account dest, int amount) {
 this.setBalance(getBalance() - amount);
 dest.setBalance(dest.getBalance() + amount);

Implementation

Goals

compute class-wide off-heap layout

 (a) remove persistent attributes
 (b) add "addr" field
 (c) generate or transform field getters/setters

 (2) replace (non-transient) field accesses
 (3) generate constructor, re-constructor
 (4) FA-wrap non-private methods

```
// transformed code (decompiled)
class Account implements PObject {
   long addr; // persistent data structure
   transient int y;
```

```
Account(String name, int balance) {
   JNVM.faStart();
   this.addr = JNVM.alloc(getClass(), size());
   this.setName(new PString(id));
   this.setBalance(balance);
   JNVM.faEnd();
```

```
•
```

. . .

```
void transferTo(Account dest, int amount) {
    JNVM.faStart();
    this.setBalance(getBalance() - amount);
    dest.setBalance(dest.getBalance() + amount);
    JNVM.faEnd();
```

Implementation

Goals

- (1) compute class-wide **off-heap layout** (a) remove persistent attributes
 - a) add "addr" field
 - (c) generate or transform field getters/setters
 - replace (non-transient) field accesses
- (3) generate constructor, re-constructor
- (4) **FA-wrap** non-private methods

Tool implementation

- (1) Bytecode-to-bytecode transformer
- (2) post-compilation Maven plugin

```
// transformed code (decompiled)
class Account implements PObject {
   long addr; // persistent data structure
   transient int y;
```

```
Account(String name, int balance) {
   JNVM.faStart();
   this.addr = JNVM.alloc(getClass(), size());
   this.setName(new PString(id));
   this.setBalance(balance);
   JNVM.faEnd();
```

```
void transferTo(Account dest, int amount) {
   JNVM.faStart();
   this.setBalance(getBalance() - amount);
   dest.setBalance(dest.getBalance() + amount);
   JNVM.faEnd();
```

. . .

Implementation

Automatic crash-consistent update usage = JNVM.faStart() *some code* JNVM.faEnd()

Per-thread persistent redo-log (inspired by Romulus)

```
Log new, free and updates
granularity = a block of PMEM
```

Do *not* log updates to "new" persistent objects (e.g. allocated within the FA-block)

J-PDT

Implementation

Persistent Data Types

- drop-in replacement for (part of) JDK e.g., string, native array, map.



Persistent vs Volatile data types (YCSB-A)

Takeaway:

- around 50% slower than volatile data types on DRAM

Implementation

- unsafe.{pwb,pfence, psync}
- NVMM block allocator
- recovery time GC (à la Makalu)
- validation = 1 bit in object header
 - makes atomic reclamation easier
 - allows deferring object liveness
 - interpreted on recovery to reclaim reachable invalid objects



List<Account> a = randAcc(100);

Bank b = new Bank(a); root.put("Bank", b); b.validate();

Implementation

- unsafe.{pwb,pfence, psync}
- NVMM block allocator
- recovery time GC (à la Makalu)
- validation = 1 bit in object header
 - makes atomic reclamation easier
 - allows deferring object liveness
 - interpreted on recovery to reclaim reachable invalid objects



List<Account> a = randAcc(100);
Bank b = new Bank(a); //Not atomic
root.put("Bank", b);
b.validate();

Implementation

- unsafe.{pwb,pfence, psync}
- NVMM block allocator
- recovery time GC (à la Makalu)
- validation = 1 bit in object header
 - makes atomic reclamation easier
 - allows deferring object liveness
 - interpreted on recovery to reclaim reachable invalid objects



List<Account> a = randAcc(100);
Bank b = new Bank(a);
root.put("Bank", b); //Not atomic
b.validate();

Implementation

- unsafe.{pwb,pfence, psync}
- NVMM block allocator
- recovery time GC (à la Makalu)
- validation = 1 bit in object header
 - makes atomic reclamation easier
 - allows deferring object liveness
 - interpreted on recovery to reclaim reachable invalid objects



List<Account> a = randAcc(100); Bank b = new Bank(a); root.put("Bank", b); b.validate();

Outline

Evaluation

(1) Introduction

- data persistence
- persistent memory
- NVMM
- why Java?
- prior art
- contribution overview
- demo

- (2) System Design of J-NVM
 - key idea
 - programming model
 - persistent objects
 - code generator
 - J-PFA
 - J-PDT
- (3) Evaluation
 - YCSB benchmark
 - recovery
- (4) Conclusion

YCSB Benchmark

Evaluation



Durable backends for Infinispan:

- PCJ = HashMap from Persistent Collections Java (JNI + PMDK)
- FS: ext4-dax

Hardware used: 4 Intel CLX 6230 HT 80-core 128GB DDR4, 4x128GB Optane (gen1)

Takeaways:

- J-NVM up to 10.5x (resp. 22.7x) than FS (resp. PCJ)
- no need for volatile cache

Recovery

Evaluation



TPC-B like benchmark 10M accounts (140 B each) client-server setting SIGKILL after 1 min

Takeaways:

- J-NVM is more than 5x faster to recover than FS
- no-need for graph traversal in some cases (e.g., only FA blocks)

Contribution

<u>Contribution</u> = J-NVM: off-heap persistent objects

Each persistent object is composed of

- a persistent data structure: unmanaged, allocated off-heap (NVMM)
- a proxy: managed, allocated on-heap (DRAM)

Pros:

- unique data representation (no data marshalling)
- recovery-time GC (not at runtime, does not scale)
- consistently faster than external designs (JNI, FS)

Cons:

- explicit free <u>but</u> common for durable data
- limited code re-use *but* safer programming model

- + automagic tool
- + library ~ no runtime changes

Conclusion

Related Work

Related Work

- Persistence History
- Early PMEM
- NVMM challenges
- NVMM Programming abstractions
- Persistence in Java
- NVMM in Java

Big summary table p.100:

Synoptic view & J-NVM positioning

1	Support		Failure-atomicity			Applicability			Efficiency					
	Model type	Language	Impl.	Granu.	Persistency	Static	FASE	PDT	Dual	Marsh.	NVM	Heap	Heap	Recovery
						types			rep.		perf.	mgmt	size	
File systems (§2.5.2)														
Ext4-DAX [2]	file system	any	user	page	imm.	no	no	no	yes	yes		N/A	N/A	
NOVA [331]	file system	any	CoW	page	imm./buff.	no	no	no	yes	yes		N/A	N/A	
Mashona [156]	library	Java (native)	append	word	imm.	no	no	no	yes	yes	-	N/A	N/A	
Persistent data type (§2.7.1)														
SOFT [353]	data type	C++	log-free	op.	d. lin.	yes	no	yes	no	no	+ + +	+ + +	+ + +	++
General-purpose constructions (§2.7.2)														
Capsules [57]	method	C++	dual-copy+CAS	word	detect. rec.	no	no	user	yes	no	++	N/A	N/A	+ + +
RECIPE [207]	method	C++	log-free	op.	imm.	no	no	user	no	no	++	N/A	N/A	+ + +
PMwCAS [314]	method+library	C++	log-free	word	imm.+linea.	no	no	user	no	no	++	N/A	N/A	+ + +
MOD [158]	method+data types	C++	shadow pag.+CAS	op.	imm.	no	no	yes	no	no	++	N/A	N/A	+ + +
NVTraverse [139]	method+library	C++	log-free	op.	d. lin.	no	no	user	no	no	++	N/A	N/A	+ + +
PRONTO [231]	library	C++	sem. log+chkpt.	op.	d. lin.	no	yes	no	yes	no	++	N/A	N/A	-
CX-PUC [103]	library	C++	2N replicas+CAS	object	d. lin.	yes	yes	no	no	no	+	++	++	+ + +
Montage [318]	library	C++	CoW+epoch reclam.	payload ¹	buff. d. lin.	yes	yes	no	no	no	+++	+ + +	+ + +	++
Mirror [140]	library	C++	log-free atomics	word	d. lin.	yes	no	yes	yes	no	+++	+ + +	+ + +	+
PREP-UC [93]	library	C++	2 replicas+ sem. log	$cache^2$	(buff.) d. lin.	yes	yes	no	yes	no	++	+ + +	++	++
Tracking [52]	method	C++	log-free	op.	detect. rec.	no	no	user	no	no	++	N/A	N/A	+ + +
PCOMB [130]	method	unspe.	CoW+CAS	object	detect. rec.	no	no	user	no	no	+++	++	+ + +	+ + +
ResPCT [192]	library	С	InCLL(undo+epoID)	epoch	buff. d. lin.	yes	yes	no	no	no	+ + + +	++	++	+ + +
Persistent transactional memory (§2.7.3)														
Romulus [102]	library	C++	dual-copy+vola. redo	word	d. lin.	yes	yes	no	no	no	+	+ + +	++	
OneFile [266]	library	C++	redo	word	d. lin.	yes	yes	no	no	no	+	++	+ + +	+ + +
Redo-PTM [103]	library	C++	N + 1 replicas+CAS	word	d. lin.	yes	yes	no	no	no	+	++	++	+ + +
Crafty [142]	library	С	undo	word	buff.	no	yes	no	no	no	N/A	++	+ + +	+ + +
Persistent heap managers (§2.7.5, §2.8.3, §2.8.4)														
Mnemosyne [309]	library+compiler	C or C++	redo	word	imm./buff.	no	yes	no	no	no	+	++	+ + +	+ + +
Atlas [83]	library+compiler	С	undo	word	buff.	no	yes	no	no	no	+	++	+ + +	+ + +
PMDK [19]	library	C & C++	undo	object	imm.	yes	yes	no	no	no	+	+	+ + +	+ + +
NVthreads [167]	library	C	CoW+redo	page	buff.	no	yes	no	no	no	N/A	++	+ + +	++
NV-heaps [92]	library	C++	undo	object	d. lin.	yes	yes	no	no	no	N/A	++	+ + +	+ + +
Makalu [63]	library (allocator)	С	undo	16B	imm.	no	no	no	no	no	+	++	+ + +	++
Ralloc [76]	library (allocator)	C++	user	8B	recoverable	no	no	no	no	no	+++	+ + +	+ + +	++
PCJ [16]	library	Java (JNI)	undo	object	imm.	yes	yes	yes	yes	yes			+ + +	+ + +
LLPL [20]	library	Java (JNI)	undo	object	imm.	yes	yes	no	yes	yes			+ + +	+ + +
MDS [298]	library	Java (JNI)	none	$cache^2$	flush-on-fail	yes	yes	no	no	yes		++	++	+ + +
Espresso [328]	compiler+runtime	Java (native)	undo	word	imm.	no	yes	no	no	no	+	+		+
AutoPersist [287]	compiler+runtime	Java (native)	undo	word	imm.	no	yes	no	no	no	-	-		+
Go-PMEM [143]	compiler+runtime	Go	undo	word	imm.	no	yes	no	no	no	+	+		+
This thesis (§4.3,	§4.5, §4.6)													
J-NVM [209]	library	Java (native)	user	user	user	yes	yes	yes	no	no	+ + +	+ + +	+ + +	+ + +
J-PDT [209]	library	Java (native)	log-free	word	imm.	yes	yes	yes	no	no	+++	+ + +	+ + +	++
J-PFA [209]	library	Java (native)	redo	block	imm.	yes	yes	yes	no	no	++	+ + +	+ + +	+ + +

Table 2.4: Comparison of various NVMM programming support libraries and systems.

¹payload: NVMM allocation payloads returned by the allocator.

²cache: the whole CPU cache hierarchy is invalidated and flushed to (persistent) memory.



- Data interoperability
- Test tools
- Killer App:
 - NVMM use for serverless apps (FaaS)

Outline

(1) Introduction

- data persistence
- persistent memory
- NVMM
- why Java?
- prior art
- contribution overview

- (2) System Design of J-NVM
 - demo
 - key idea
 - programming model
 - persistent objects
 - code generator
 - J-PFA
 - J-PDT
- (3) Evaluation
 - YCSB benchmark
 - recovery

(4) Conclusion