Neverlast and Beyond

Handling Volatility in NVRAM-centric Operating Systems

GI Fachgruppe Betriebssysteme - Herbsttreffen September 21, 2021

Christian Eichler¹, Henriette Hofmeier¹, Wolfgang Schröder-Preikschat², and Timo Hönig¹

¹Ruhr-Universität Bochum (RUB) ²Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU)





Energy-harvesting Systems...

• Obtain power from the environment (e.g., solar panels)





Energy-harvesting Systems...

- Obtain power from the environment (e.g., solar panels)
- Depend on environmental conditions (e.g., sun, clouds)





Energy-harvesting Systems...

- Obtain power from the environment (e.g., solar panels)
- Depend on environmental conditions (e.g., sun, clouds)
- Often come with limited energy storage (e.g., capacitors)





Energy-harvesting Systems...

- Obtain power from the environment (e.g., solar panels)
- Depend on environmental conditions (e.g., sun, clouds)
- Often come with limited energy storage (e.g., capacitors)
- May lose power supply anytime





Energy-harvesting Systems...

- Obtain power from the environment (e.g., solar panels)
- Depend on environmental conditions (e.g., sun, clouds)
- Often come with limited energy storage (e.g., capacitors)
- May lose power supply anytime



→ Intermittently-powered systems may run out of power, but **must not lose system state**



Non-volatile memory retains data on power outage

- Ideal solution for energyharvesting applications
- + Can be used as drop-in replacement for conventional volatile main memory

Non-volatile memory retains data on power outage 🄰

- Ideal solution for energyharvesting applications
- + Can be used as drop-in replacement for conventional volatile main memory
- + Commercially available







Non-volatile memory retains data on power outage

- Ideal solution for energyharvesting applications
- + Can be used as drop-in replacement for conventional volatile main memory
- + Commercially available

But some volatile components remain:

- CPU registers & caches
- Hardware configuration







Non-volatile memory retains data on power outage 🎴

- Ideal solution for energyharvesting applications
- + Can be used as drop-in replacement for conventional volatile main memory
- + Commercially available

But some volatile components remain:

- CPU registers & caches
- Hardware coi

 \rightarrow Data in volatile components must be taken into account when persisting the system state



```
int sensor = open(TEMPERATURE_SENSOR);
```



```
int sensor = open(TEMPERATURE_SENSOR);
write(sensor, "oversampling x16");
```



```
int sensor = open(TEMPERATURE_SENSOR);
write(sensor, "oversampling x16");
power failure
```



```
int sensor = open(TEMPERATURE_SENSOR);
write(sensor, "oversampling x16");
power failure
```

Bootup after power failure



```
int sensor = open(TEMPERATURE_SENSOR);
write(sensor, "oversampling x16");
powerfailure
```

```
Bootup after power failure
write(sensor, "offset +2");
```



```
int sensor = open(TEMPERATURE_SENSOR);
write(sensor, "oversampling x16");
powerfailure
```

Bootup after power failure
write(sensor, "offset +2");



Expected Configuration
Oversampling: x16
Offset: +2

```
int sensor = open(TEMPERATURE_SENSOR);
write(sensor, "oversampling x16");
powerfailure
```

Bootup after power failure
write(sensor, "offset +2");



Expected Configuration

Oversampling: x16

\rightarrow Losing volatile state eventually leads to inconsistencies



Design & Implementation of Neverlast

Neverlast: Handling Device and System Consistency

Consistency comprises...

- Application and OS data
- Volatile cache and register contents
- But also volatile hardware configuration

Therefore, Neverlast...

- 1. Minimizes the volatile state by running in NVRAM
- 2. Saves registers & caches on power outage
- 3. Traces device-configuration syscalls



Neverlast: Detecting & Handling Power Outages[1]

Detecting Power Outages

1. Voltage comparator monitors V_{CC} and interrupts the CPU when V_{CC} falls below threshold





 [1] Neverlast: Towards the Design and Implementation of the NVM-based Everlasting Operating System. Christian Eichler, Henriette Hofmeier, Stefan Reif, Timo Hönig, Jörg Nolte, and Wolfgang Schröder-Preikschat. In: Proceedings of the 54th Hawai'i International Conference on System Sciences (HICSS-54)

Neverlast: Detecting & Handling Power Outages[1]

Detecting Power Outages

- 1. Voltage comparator monitors V_{CC} and interrupts the CPU when V_{CC} falls below threshold
- 2. ISR stores remaining volatile data (registers)
- 3. Marks stored system state as valid
- 4. Stops system execution





 Neverlast: Towards the Design and Implementation of the NVM-based Everlasting Operating System. Christian Eichler, Henriette Hofmeier, Stefan Reif, Timo Hönig, Jörg Nolte, and Wolfgang Schröder-Preikschat. In: Proceedings of the 54th Hawai'i International Conference on System Sciences (HICSS-54)

Neverlast: Detecting & Handling Power Outages[1]

Detecting Power Outages

- 1. Voltage comparator monitors V_{CC} and interrupts the CPU when V_{CC} falls below threshold
- 2. ISR stores remaining volatile data (registers)
- 3. Marks stored system state as valid
- 4. Stops system execution

Recovering from Power Outages

- 1. Check stored system state for validity
- 2. Restore volatile state
- 3. Continue execution

 Neverlast: Towards the Design and Implementation of the NVM-based Everlasting Operating System. Christian Eichler, Henriette Hofmeier, Stefan Reif, Timo Hönig, Jörg Nolte, and Wolfgang Schröder-Preikschat. In: Proceedings of the 54th Hawai'i International Conference on System Sciences (HICSS-54)







Interrupt-based Mechanism Problematic for Peripheral Devices

- Unknown device state
- Configuration not readable

- Slow communication bus
- High power consumption

Interrupt-based Mechanism Problematic for Peripheral Devices

- Unknown device state
- Configuration not readable

- Slow communication bus
- High power consumption

Neverlast's Approach: Trace System Calls

- Transaction-like behavior (via replay_start and replay_stop)
- Log system calls (including parameters) to persistent log
- Replay log after power outage

[2] Neverlast: an NVM-centric Operating System for Persistent Edge Systems. Christian Eichler, Henriette Hofmeier, Stefan Reif, Timo Hönig, Jörg Nolte, and Wolfgang Schröder-Preikschat. In: Proceedings of the 12th ACM SIGOPS Asia-Pacific Workshop on Systems (APSys 2021)

Interrupt-based Mechanism Problematic for Peripheral Devices

- Unknown device state
- Configuration not readable

- Slow communication bus
- High power consumption

Neverlast's Approach: Trace System Calls

- Transaction-like behavior (via replay_start and replay_stop)
- Log system calls (including parameters) to persistent log
- Replay log after power outage

\rightarrow Use **syscall-based transaction mechanism** for peripheral devices

[2] Neverlast: an NVM-centric Operating System for Persistent Edge Systems. Christian Eichler, Henriette Hofmeier, Stefan Reif, Timo Hönig, Jörg Nolte, and Wolfgang Schröder-Preikschat. In: Proceedings of the 12th ACM SIGOPS Asia-Pacific Workshop on Systems (APSys 2021)

```
1void app_main() {
```

- 2 int sen = open(TEMP_SENSOR)
- 3 replay_start(sen, 1)
- 4 write(sen, "oversampling x16")

```
5 write(sen, "offset +2")
6 replay_stop(sen)
7 /*... */
8}
```

```
1void app main() {
 int sen = open(TEMP SENSOR)
2
 replay start(sen, 1)
3
4 write(sen, "oversampling x16")
          power-failure interrupt -
5 write(sen, "offset +2")
6 replay stop(sen)
 /* . . . */
8 }
```

Sensor Configuration

Oversampling: x16 Offset: +0



```
1void app main() {
 int sen = open(TEMP SENSOR)
2
 replay start(sen, 1)
3
4 write(sen, "oversampling x16")
          power-failure interrupt -
5 write(sen, "offset +2")
6 replay stop(sen)
 /* . . . */
8 }
```

Sensor Configuration

Oversampling: - Offset: -

```
//store volatile registers
pushm.a #12, r15
poweroff()
```

```
1void app main() {
2 int sen = open(TEMP SENSOR)
 replay start(sen, 1)
3
4 write(sen, "oversampling x16")
         power-failure interrupt -
5 write(sen, "offset +2")
6 replay stop(sen)
7 /*... */
8 }
```

Sensor Configuration

Oversampling: x16 Offset: +0

```
//store volatile registers
pushm.a #12, r15
poweroff()
```

//on power recovery
int sen = open(TEMP_SENSOR)
write(sen, "oversampling x16")
popm.a #12, r15
return

```
1void app main() {
2 int sen = open(TEMP SENSOR)
 replay start(sen, 1)
3
4 write(sen, "oversampling x16")
         power-failure interrupt
5 write(sen, "offset +2")
6 replay stop(sen)
 /* . . . */
8 }
```

Sensor Configuration

Oversampling: x16 Offset: +2

```
//store volatile registers
pushm.a #12, r15
poweroff()
```

//on power recovery
int sen = open(TEMP_SENSOR)
write(sen, "oversampling x16")
popm.a #12, r15
return

Evaluation

MSP-EXP430FR5994 LaunchPad Development Kit

- 16 bit RISC processor, up to 24 MHz
- 8 KB of SRAM
- 256 KB of FRAM with 2-way-2-line read cache (2x 8 Byte)
- Remote-controlled via relais card
- Shunt-based energy measurement Tektronix MSO 4034



Evaluation - Power Consumption by Supply Monitoring



Evaluation - Power Consumption by Supply Monitoring



Evaluation - Power Consumption by Supply Monitoring



Evaluation - Performance of NVRAM vs. RAM (sequential read/write)



Evaluation - Performance of NVRAM vs. RAM (non-sequential read/write)



Evaluation - Performance of NVRAM vs. RAM (non-sequential read/write)



- ? Energy-harvesting systems may lose power anytime
- ? Consistency in systems with mixed volatility...
- ? and peripheral devices





- ✓ Energy-harvesting systems may lose power anytime → Monitor power supply & detect power outages
- ? Consistency in systems with mixed volatility...
- ? and peripheral devices





- ✓ Energy-harvesting systems may lose power anytime → Monitor power supply & detect power outages
- Consistency in systems with mixed volatility...
 - ightarrow Save volatile system state (registers) on power failure
- ? and peripheral devices





- ✓ Energy-harvesting systems may lose power anytime → Monitor power supply & detect power outages
- Consistency in systems with mixed volatility...
 - \rightarrow Save volatile system state (registers) on power failure
- ✓ and peripheral devices
 - \rightarrow Trace & replay device-configuration syscalls





A Peak into the Future...





Persistent memory used to be (really) slow and high latency

- Think of: Mechanical hard drives
- OSes avoid persistent memory when possible
- ightarrow Use of volatile caches



Persistent memory used to be (really) slow and high latency

- Think of: Mechanical hard drives
- OSes avoid persistent memory when possible
- ightarrow Use of volatile caches



The OS handles RAM volatility in software:

- Periodic flush of write-back file caches
- Periodic flush of the superblock & inode cache
- "Hibernate to disk"

Imagine using (only) non-volatile main memory:



Imagine using (only) non-volatile main memory:

Previously volatile data is no longer volatile
 OS persistency measures less/no longer important



Imagine using (only) non-volatile main memory:

- Previously volatile data is no longer volatile
 → OS persistency measures less/no longer important
- But volatile hardware components remain: caches, registers, device configuration



Imagine using (only) non-volatile main memory:

- Previously volatile data is no longer volatile
 OS persistency measures less/no longer important
- But volatile hardware components remain: caches, registers, device configuration
- **DFG** NEON / Linux on Intel Optane:
 - Port Neverlast's ideas to Linux
 - Identify & remove obsolete persistency measures



NVRAM-centric operating systems come with **challenges and opportunities**:

- ? Persistency & consistency in presence of volatile components
- + Optimization of the OS:
 - + Simplifications
 - Reduced TCB
 - + Reduced jitter





Bundesministerium für Bildung und Forschung 16KIS1315



Conclusion

NVRAM-centric operating systems come with **challenges and opportunities**:

- ? Persistency & consistency in presence of volatile components
- + Optimization of the OS:
 - + Simplifications
 - Reduced TCB
 - + Reduced jitter





Bundesministerium für Bildung und Forschung 16KIS1315

