# Efficient Guest Agnostic Virtualization With Embedded Power Architecture® KVM
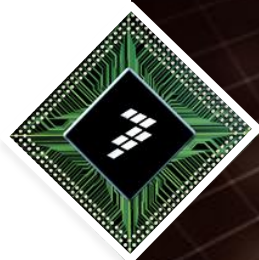
Varun Sethi
Freescale Semiconductor

# Motivation

- Growing interest in virtualization using KVM on embedded Power Architecture platforms

- Requirement to run various customer specific operating systems with embedded Power Architecture KVM

- Possible to run unmodified guest on embedded Power Architecture KVM – But this comes at a significant cost associated with VM exits

  - *Problem severe for cores without virtualization assists*

- Possible to "Binary translate" guest privileged instructions from the host side.

  - Continue to run an unmodified guest

# Adaptive Binary Translation

- Dynamically infer the instructions that cause a large number of VM exits

- Binary Translate these instructions to a faster emulation code .

- Instruction translation maintained on a shared memory region

- "Memory tracing" (remove R/W permissions to shared page) implemented to control guest access to modified pages

  - Guest access to page generates a DSI exception



Guest OS

VM Exit

Shared Pages

Host Kernel    KVM    LABT

1. Find hot instructions
2. Adaptive Patching
3. Tracing

# Design Challenges

- Handling complex instruction translations
    - Multi line patching could be complex
    - Translation cache placement issue

- Minimizing performance overhead (excessive DSIs) due to Memory tracing
    - Use of huge pages (TLB1)
    - Self referential code
        - Access to sys call table and exception prolog in case of Linux
    - Self modifying code
        - Code from the modified page trying to modify (Write) code on the same page

*freescale* ™

# Solutions to Design Challenges



- Adaptive data mirroring algorithm to address self referential code

  ▪ Data causing excessive DSIs mirrored to a guest memory region having R/W permission
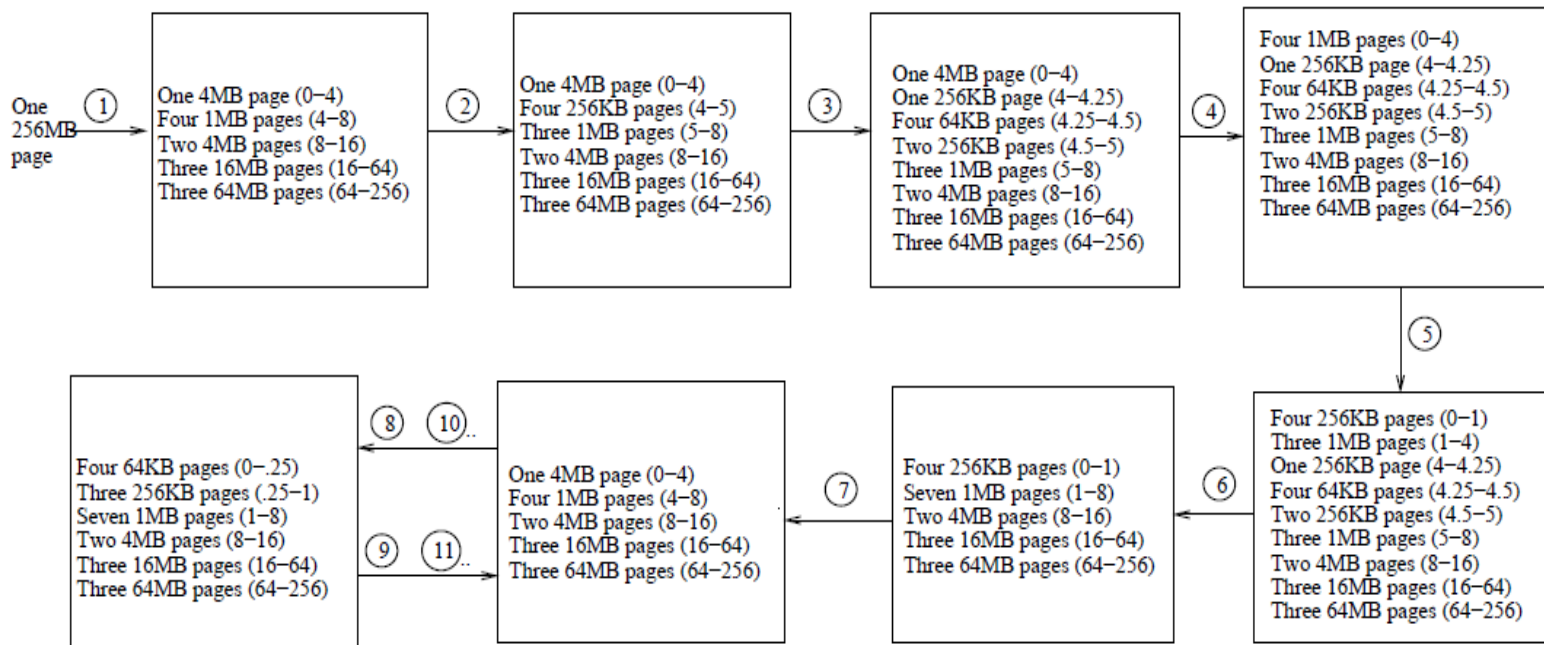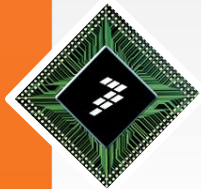
- Addressing Translation Cache placement issue

  – Mechanism developed for stealing space from read only guest section

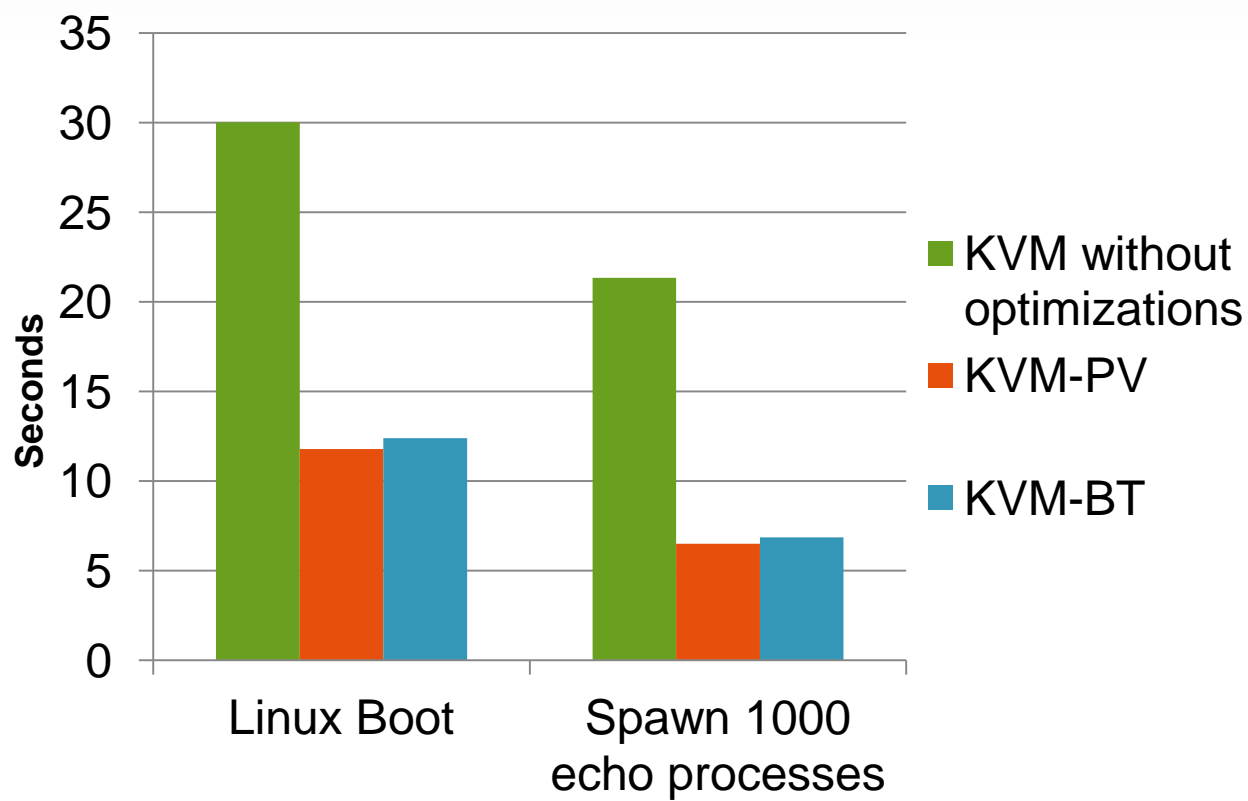  – The data from the stolen section mirrored to new location

One 256MB page

① One 4MB page (0–4)
Four 1MB pages (4–8)
Two 4MB pages (8–16)
Three 16MB pages (16–64)
Three 64MB pages (64–256)

② One 4MB page (0–4)
Four 256KB pages (4–5)
Three 1MB pages (5–8)
Two 4MB pages (8–16)
Three 16MB pages (16–64)
Three 64MB pages (64–256)

③ One 4MB page (0–4)
One 256KB page (4–4.25)
Four 64KB pages (4.25–4.5)
Two 256KB pages (4.5–5)
Three 1MB pages (5–8)
Two 4MB pages (8–16)
Three 16MB pages (16–64)
Three 64MB pages (64–256)

④ Four 1MB pages (0–4)
One 256KB page (4–4.25)
Four 64KB pages (4.25–4.5)
Two 256KB pages (4.5–5)
Three 1MB pages (5–8)
Two 4MB pages (8–16)
Three 16MB pages (16–64)
Three 64MB pages (64–256)

⑤ Four 256KB pages (0–1)
Three 1MB pages (1–4)
One 256KB page (4–4.25)
Four 64KB pages (4.25–4.5)
Two 256KB pages (4.5–5)
Three 1MB pages (5–8)
Two 4MB pages (8–16)
Three 16MB pages (16–64)
Three 64MB pages (64–256)

⑥ Four 256KB pages (0–1)
Seven 1MB pages (1–8)
Two 4MB pages (8–16)
Three 16MB pages (16–64)
Three 64MB pages (64–256)

⑦ One 4MB page (0–4)
Four 1MB pages (4–8)
Two 4MB pages (8–16)
Three 16MB pages (16–64)
Three 64MB pages (64–256)

⑧ ⑩ ⑨ ⑪ Four 64KB pages (0–.25)
Three 256KB pages (.25–1)
Seven 1MB pages (1–8)
Two 4MB pages (8–16)
Three 16MB pages (16–64)
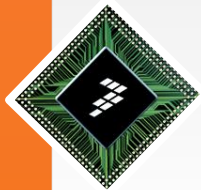Three 64MB pages (64–256)

- Addressing Memory tracing issues:
  - Adaptive page resizing algorithm for addressing issue arising out of the use of Huge TLB and self modifying code
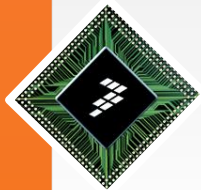    - Dynamically splitting/merging traced pages

# Performance

# Summary

- With "Adaptive Binary Translation" it's possible to run an unmodified guest  efficiently with embedded Power Architecture KVM

- Possible to mitigate memory tracing overhead using "Adaptive Page Resizing"  and  "Adaptive Data Mirroring"

- Performance is comparable to the existing PV solution

**freescale** ™

# Questions?

**freescale** ™

freescale™