

Echtzeitfähigkeit virtueller Maschinen

Robert Kaiser¹ Dieter Zöbel²

¹Fachbereich 2 – Technik, Informatik und Wirtschaft
Fachhochschule Bingen

²Fachbereich 4 – Informatik
Universität Koblenz-Landau

Echtzeit 2010, Boppard, 18.11.2008



1 Einleitung

- Motivation
- Herausforderungen

2 Planung unter Virtualisierungsumgebungen

- Hierarchischer Scheduler
- Planbarkeit

3 Kosten der Virtualisierung

- Planungsbedingte Verluste
- Umschaltkosten

4 Zusammenfassung

1 Einleitung

- Motivation
- Herausforderungen

2 Planung unter Virtualisierungsumgebungen

- Hierarchischer Scheduler
- Planbarkeit

3 Kosten der Virtualisierung

- Planungsbedingte Verluste
- Umschaltkosten

4 Zusammenfassung

Rechenleistung eingebetteter Systeme steigt überproportional gegenüber Kosten

- Potenzial ...
 - ... **wirtschaftlich**: höhere Rechenleistung bei geringeren Kosten
 - ... **technisch**: weniger Bauteile ⇒ geringere Ausfallwahrscheinlichkeit
- ⇒ Ziel: Mehr Funktionen auf einem Rechner zusammenfassen
- ⇒ Insgesamt weniger, dafür leistungsfähigere Rechner

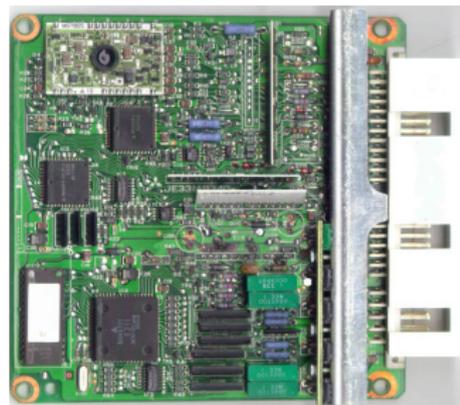


Abbildung: KFZ-Steuergerät

⇒ Konsolidierung eingebetteter Systeme steht an

Rechenleistung eingebetteter Systeme steigt überproportional gegenüber Kosten

- Potenzial ...
 - ... **wirtschaftlich**: höhere Rechenleistung bei geringeren Kosten
 - ... **technisch**: weniger Bauteile ⇒ geringere Ausfallwahrscheinlichkeit
- ⇒ Ziel: Mehr Funktionen auf einem Rechner zusammenfassen
- ⇒ Insgesamt weniger, dafür leistungsfähigere Rechner

⇒ Konsolidierung eingebetteter Systeme steht an



Abbildung: Verteilung von Steuergeräten in einem PKW

Rechenleistung eingebetteter Systeme steigt überproportional gegenüber Kosten

- Potenzial ...
 - ... **wirtschaftlich**: höhere Rechenleistung bei geringeren Kosten
 - ... **technisch**: weniger Bauteile ⇒ geringere Ausfallwahrscheinlichkeit
- ⇒ Ziel: Mehr Funktionen auf einem Rechner zusammenfassen
- ⇒ Insgesamt weniger, dafür leistungsfähigere Rechner

⇒ Konsolidierung eingebetteter Systeme steht an



Abbildung: Verteilung von Steuergeräten in einem PKW

Rechenleistung eingebetteter Systeme steigt überproportional gegenüber Kosten

- Potenzial ...
 - ... **wirtschaftlich**: höhere Rechenleistung bei geringeren Kosten
 - ... **technisch**: weniger Bauteile ⇒ geringere Ausfallwahrscheinlichkeit
- ⇒ Ziel: Mehr Funktionen auf einem Rechner zusammenfassen
- ⇒ Insgesamt weniger, dafür leistungsfähigere Rechner

⇒ Konsolidierung eingebetteter Systeme steht an



Abbildung: Verteilung von Steuergeräten in einem PKW

Herausforderungen:

- Ausschließen ungewollter Wechselwirkungen
(„natürliche“ Trennung der Subsysteme aufrecht erhalten)
- Verschiedene Laufzeitumgebungen auf einem Rechnerknoten
(Jeder Anwendung ihre „eigene Welt“)

Ähnlich den Anforderungen bei der Serverkonsolidierung

- Mehrere virtuelle Maschinen auf einem physischen Rechnerknoten
- Gegenseitige Isolierung, Aufteilung der Betriebsmittel

Aber: ist Virtualisierung „Echtzeitfähig“?

Herausforderungen:

- Ausschließen ungewollter Wechselwirkungen
(„natürliche“ Trennung der Subsysteme aufrecht erhalten)
- Verschiedene Laufzeitumgebungen auf einem Rechnerknoten
(Jeder Anwendung ihre „eigene Welt“)

Ähnlich den Anforderungen bei der Serverkonsolidierung

- Mehrere virtuelle Maschinen auf einem physischen Rechnerknoten
- Gegenseitige Isolierung, Aufteilung der Betriebsmittel

Aber: ist Virtualisierung „Echtzeitfähig“?

1 Einleitung

- Motivation
- Herausforderungen

2 Planung unter Virtualisierungsumgebungen

- Hierarchischer Scheduler
- Planbarkeit

3 Kosten der Virtualisierung

- Planungsbedingte Verluste
- Umschaltkosten

4 Zusammenfassung

Bei bisherigen Virtualisierungsansätzen:

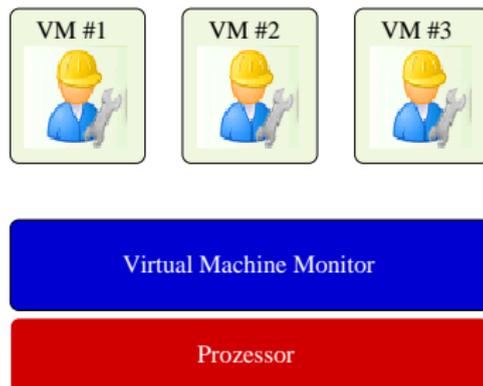
- Einfluss der Virtualisierung auf das Zeitverhalten wurde nicht betrachtet (vgl. [Popek/Goldberg])

Für eingebettete Systeme nicht akzeptabel

- Prozessor muss im Zeitmultiplex-Verfahren zugeteilt werden
- ⇒ VMM agiert als *globaler* Scheduler
- Virtuelle Maschinen enthalten i.A. Betriebssysteme mit Anwendungen

⇒ Betriebssysteme agieren als *lokale* Scheduler

Beide Scheduler-Ebenen sind zu betrachten



Bei bisherigen Virtualisierungsansätzen:

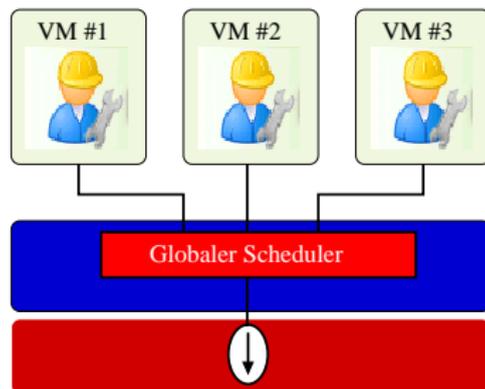
- Einfluss der Virtualisierung auf das Zeitverhalten wurde nicht betrachtet (vgl. [Popek/Goldberg])

Für eingebettete Systeme nicht akzeptabel

- Prozessor muss im Zeitmultiplex-Verfahren zugeteilt werden
- ⇒ VMM agiert als *globaler* Scheduler
- Virtuelle Maschinen enthalten i.A. Betriebssysteme mit Anwendungen

⇒ Betriebssysteme agieren als *lokale* Scheduler

Beide Scheduler-Ebenen sind zu betrachten



Bei bisherigen Virtualisierungsansätzen:

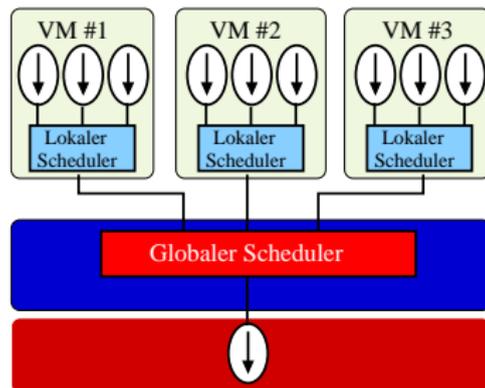
- Einfluss der Virtualisierung auf das Zeitverhalten wurde nicht betrachtet (vgl. [Popek/Goldberg])

Für eingebettete Systeme nicht akzeptabel

- Prozessor muss im Zeitmultiplex-Verfahren zugeteilt werden
- ⇒ VMM agiert als *globaler* Scheduler
- Virtuelle Maschinen enthalten i.A. Betriebssysteme mit Anwendungen

⇒ Betriebssysteme agieren als *lokale* Scheduler

Beide Scheduler-Ebenen sind zu betrachten



Bei bisherigen Virtualisierungsansätzen:

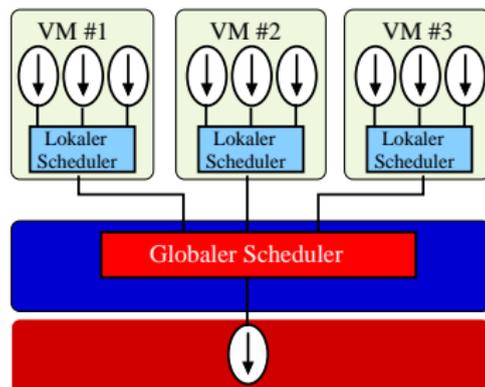
- Einfluss der Virtualisierung auf das Zeitverhalten wurde nicht betrachtet (vgl. [Popek/Goldberg])

Für eingebettete Systeme nicht akzeptabel

- Prozessor muss im Zeitmultiplex-Verfahren zugeteilt werden
- ⇒ VMM agiert als *globaler* Scheduler
- Virtuelle Maschinen enthalten i.A. Betriebssysteme mit Anwendungen

⇒ Betriebssysteme agieren als *lokale* Scheduler

Beide Scheduler-Ebenen sind zu betrachten



Bei bisherigen Virtualisierungsansätzen:

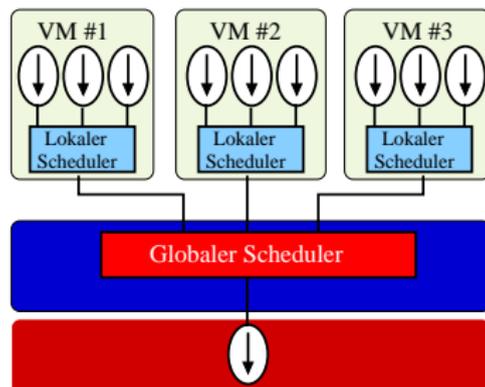
- Einfluss der Virtualisierung auf das Zeitverhalten wurde nicht betrachtet (vgl. [Popek/Goldberg])

Für eingebettete Systeme nicht akzeptabel

- Prozessor muss im Zeitmultiplex-Verfahren zugeteilt werden
- ⇒ VMM agiert als *globaler* Scheduler
- Virtuelle Maschinen enthalten i.A. Betriebssysteme mit Anwendungen

⇒ Betriebssysteme agieren als *lokale* Scheduler

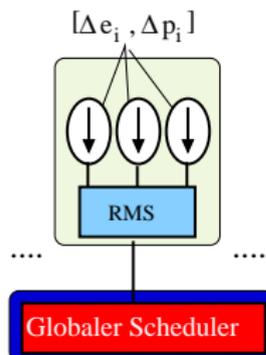
Beide Scheduler-Ebenen sind zu betrachten



Prozessmodell nach Liu/Layland:

• Gegeben:

- 1 $P = \{1, \dots, n\}$: Menge period. Prozesse mit:
 - Ausführungszeiten Δe_i
 - Periodendauern Δp_i
- 2 Lokaler Scheduler (hier: RMS)

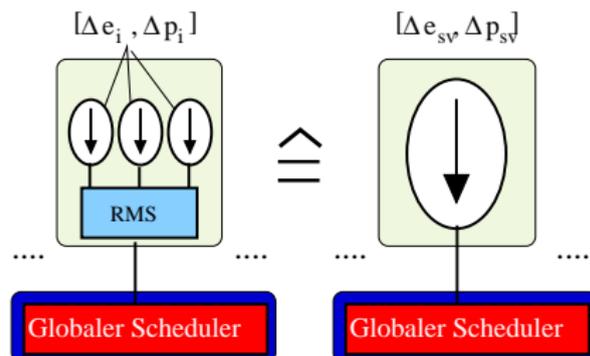


- **Gesucht:** Äquivalente Darstellung der Prozessmenge als ein periodischer Stellvertreterprozess mit $\Delta e_{SV}, \Delta p_{SV}$.

Prozessmodell nach Liu/Layland:

- **Gegeben:**

- 1 $P = \{1, \dots, n\}$: Menge period. Prozesse mit:
 - Ausführungszeiten Δe_i
 - Periodendauern Δp_i
- 2 Lokaler Scheduler (hier: RMS)



- **Gesucht:** Äquivalente Darstellung der Prozessmenge als ein periodischer Stellvertreterprozess mit $\Delta e_{sv}, \Delta p_{sv}$.

Planen nach Monotonen Raten

- Vorschrift: Prioritäten umgekehrt proportional zu Periodendauern festlegen:

$$\Delta p_i > \Delta p_j \Leftrightarrow \text{prio}(i) < \text{prio}(j)$$

- Ansatz für Stellvertreterprozess:

- 1 Fasse alle **nicht** zugeteilten Zeiten als periodischen „Störprozess“ δ auf.
- 2 Wenn $\Delta p_{sv} \leq \min(\Delta p_i)$ kann δ als neuer höchstpriorer Prozess zu P hinzugefügt werden.

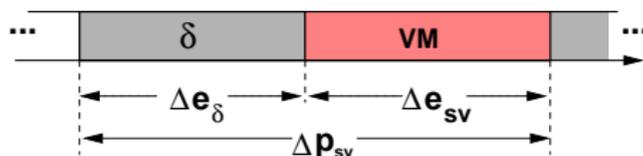
⇒ **RMS bleibt anwendbar**

Planen nach Monotonen Raten

- Vorschrift: Prioritäten umgekehrt proportional zu Periodendauern festlegen:

$$\Delta p_i > \Delta p_j \Leftrightarrow \text{prio}(i) < \text{prio}(j)$$

- Ansatz für Stellvertreterprozess:



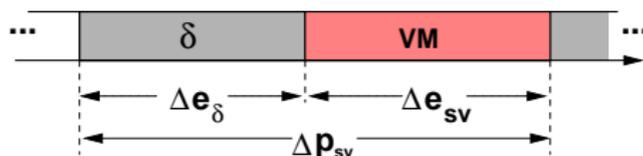
- 1 Fasse alle **nicht** zugeteilten Zeiten als periodischen „Störprozess“ δ auf.
- 2 Wenn $\Delta p_{sv} \leq \min(\Delta p_i)$ kann δ als neuer höchstpriorer Prozess zu P hinzugefügt werden.

⇒ RMS bleibt anwendbar

- Vorschrift: Prioritäten umgekehrt proportional zu Periodendauern festlegen:

$$\Delta p_i > \Delta p_j \Leftrightarrow \text{prio}(i) < \text{prio}(j)$$

- Ansatz für Stellvertreterprozess:



- 1 Fasse alle **nicht** zugeteilten Zeiten als periodischen „Störprozess“ δ auf.
- 2 Wenn $\Delta p_{sv} \leq \min(\Delta p_i)$ kann δ als neuer höchstpriorer Prozess zu P hinzugefügt werden.

⇒ **RMS bleibt anwendbar**

Frage: Bis zu welcher Auslastung ist der Plan brauchbar?

- Abgewandelte Form des Beweises von Liu/Layland (→ [Kaiser/Zöbel ETFA09]). Ergebnis:

$$U(P)_{min} = n \cdot \left(\sqrt[n]{\frac{2}{U_\delta + 1}} - 1 \right)$$
$$\lim_{n \rightarrow \infty} U(P)_{min} = \ln \left(\frac{2}{U_\delta + 1} \right)$$

- Dabei ist: $U_\delta = \frac{\Delta e_\delta}{\Delta p_{sv}}$ (Auslastung durch den Störprozess δ)

- Periodendauer lt. Voraussetzung (s.o.): $\Delta p_{sv} = \min(\Delta p_i)$
- Außerdem gilt (s.o.): $\Delta e_{sv} + \Delta e_{\delta} = \Delta p_{sv}$
 $\Rightarrow U_{sv} = \frac{\Delta e_{sv}}{\Delta p_{sv}} = 1 - U_{\delta}$

$$\Rightarrow \Delta e_{sv} = 2 \cdot \Delta p_{sv} \cdot \left(1 - \frac{1}{\left(\frac{U(P)}{n} + 1 \right)^n} \right)$$
$$\lim_{n \rightarrow \infty} \Rightarrow \Delta e_{sv} = 2 \cdot \Delta p_{sv} \cdot \left(1 - e^{-U(P)} \right)$$

- \Rightarrow Damit: $\Delta e_{sv}, \Delta p_{sv}$ für alle VMs berechenbar
- \Rightarrow Globaler Scheduler plant wiederum periodische Prozesse
- \Rightarrow Kann seinerseits z.B. RMS oder EDF einsetzen

1 Einleitung

- Motivation
- Herausforderungen

2 Planung unter Virtualisierungsumgebungen

- Hierarchischer Scheduler
- Planbarkeit

3 Kosten der Virtualisierung

- Planungsbedingte Verluste
- Umschaltkosten

4 Zusammenfassung

Verluste sind unvermeidlich:

- Es gibt Zeitpunkte, zu denen eine VM (d.h. ihr Stellvertreterprozess) den Prozessor besitzt, aber kein lokaler Prozess ist rechenwillig
- D.h. mit:
 - $U_{sv} = \frac{\Delta e_{sv}}{\Delta p_{sv}}$ Äußere Auslastung
 - $U(P) = \sum_{i \in P} \frac{\Delta e_i}{\Delta p_i}$ Innere Auslastung durch P
- gilt: *relative Auslastung*: $u_{sv} = \frac{U_{sv}}{U(P)} > 1$

Maß für planungsbedingte Verluste (für große Prozessanzahlen):

$$u_{sv} = \frac{2}{U(P)} \cdot (1 - e^{-U(P)})$$

Für kleine Einzel-Auslastungen $U(P)$ gilt $u_{sv} \rightarrow 2$

⇒ Planungsbedingte Verluste sind beschränkt (auf maximal 50%)

Verluste sind unvermeidlich:

- Es gibt Zeitpunkte, zu denen eine VM (d.h. ihr Stellvertreterprozess) den Prozessor besitzt, aber kein lokaler Prozess ist rechenwillig
- D.h. mit:
 - $U_{sv} = \frac{\Delta e_{sv}}{\Delta p_{sv}}$ Äußere Auslastung
 - $U(P) = \sum_{i \in P} \frac{\Delta e_i}{\Delta p_i}$ Innere Auslastung durch P
- gilt: *relative Auslastung*: $u_{sv} = \frac{U_{sv}}{U(P)} > 1$

Maß für planungsbedingte Verluste (für große Prozessanzahlen):

$$u_{sv} = \frac{2}{U(P)} \cdot (1 - e^{-U(P)})$$

Für kleine Einzel-Auslastungen $U(P)$ gilt $u_{sv} \rightarrow 2$

⇒ Planungsbedingte Verluste sind beschränkt (auf maximal 50%)

In bisherigen Betrachtungen wurde Prozesswechsel-Aufwand (wie oft) vernachlässigt, aber:

- Bedingung für Anwendbarkeit RMS (s.o.) kann zu hohen Umschaltfrequenzen führen
- „Dynamische Umnutzung“ planungsbedingter Verluste prinzipiell möglich, aber je nach Umschaltfrequenz nicht sinnvoll

⇒ **Quantitatives Erfassen der Umschaltkosten ist nötig**

- Problem: Wg. Caches ist eine analytische Beschreibung dieser Kosten i.d.R. nicht möglich
 - Volle Kenntnis über Zustand und Verhalten der VMs wäre erforderlich
 - Selbst dann: Berechnung für on-line Scheduler zu komplex
- Ansatz hier: Empirische Vorgehensweise
 - Beobachten/Messen des Systemverhaltens
 - Approximation durch mit den Messwerten parametrisierte, funktionale Beschreibung

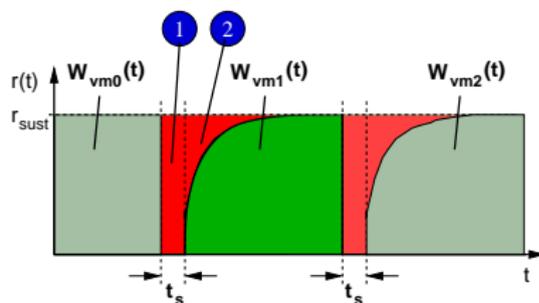
Umschaltverluste zusammengesetzt aus:

- 1 Unterbrechungszeit (fix)
- 2 Cache-bedingten Kosten (Zeitabhängig):

- Einfaches Modell des Cache-Verhaltens
- Worst-Case Verhalten ist reproduzierbar (und messbar)

⇒ Vorgehensweise:

- Worst-Case Verhalten messen
- Geeignete Funktionen mit Messwerten parametrieren
- Verluste anhand Scheduler-Simulation aufsummieren

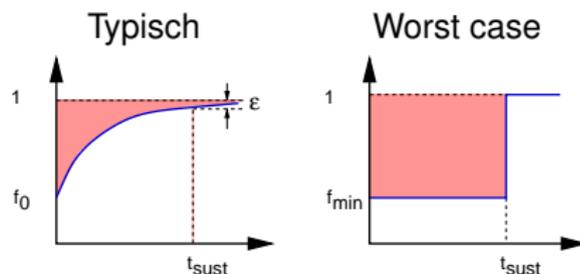


Umschaltverluste zusammengesetzt aus:

- 1 Unterbrechungszeit (fix)
- 2 Cache-bedingten Kosten (Zeitabhängig):
 - Einfaches Modell des Cache-Verhaltens
 - Worst-Case Verhalten ist reproduzierbar (und messbar)

⇒ **Vorgehensweise:**

- Worst-Case Verhalten messen
- Geeignete Funktionen mit Messwerten parametrieren
- Verluste anhand Scheduler-Simulation aufsummieren

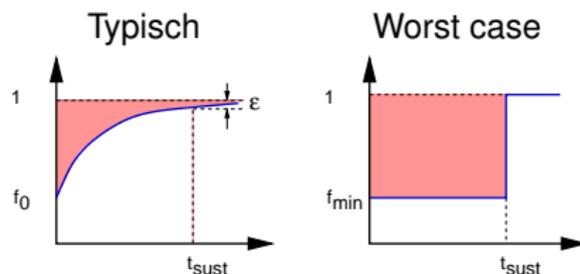


Umschaltverluste zusammengesetzt aus:

- 1 Unterbrechungszeit (fix)
- 2 Cache-bedingten Kosten (Zeitabhängig):
 - Einfaches Modell des Cache-Verhaltens
 - Worst-Case Verhalten ist reproduzierbar (und messbar)

⇒ Vorgehensweise:

- Worst-Case Verhalten messen
- Geeignete Funktionen mit Messwerten parametrieren
- Verluste anhand Scheduler-Simulation aufsummieren



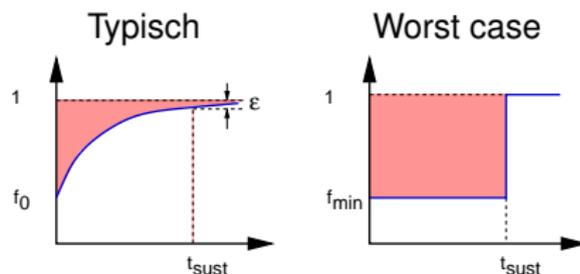
Umschaltverluste zusammengesetzt aus:

- 1 Unterbrechungszeit (fix)
- 2 Cache-bedingten Kosten (Zeitabhängig):

- Einfaches Modell des Cache-Verhaltens
- Worst-Case Verhalten ist reproduzierbar (und messbar)

⇒ Vorgehensweise:

- Worst-Case Verhalten messen
- Geeignete Funktionen mit Messwerten parametrieren
- Verluste anhand Scheduler-Simulation aufsummieren

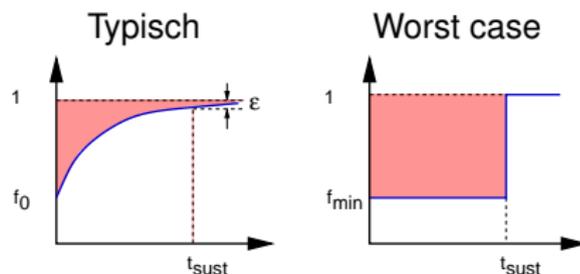


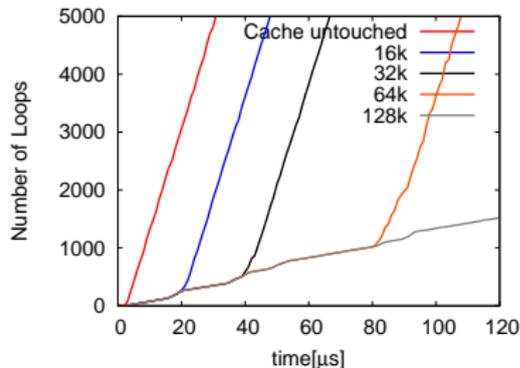
Umschaltverluste zusammengesetzt aus:

- 1 Unterbrechungszeit (fix)
- 2 Cache-bedingten Kosten (Zeitabhängig):
 - Einfaches Modell des Cache-Verhaltens
 - Worst-Case Verhalten ist reproduzierbar (und messbar)

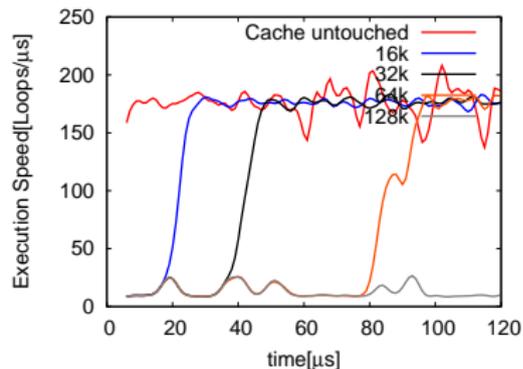
⇒ Vorgehensweise:

- Worst-Case Verhalten messen
- Geeignete Funktionen mit Messwerten parametrieren
- Verluste anhand Scheduler-Simulation aufsummieren





- Gemessen: Arbeit
(Schleifendurchläufe pro Zeiteinheit,
verschiedene „working sets“)



- Ableitung: Fortschrittsrate
(= Rechenleistung)

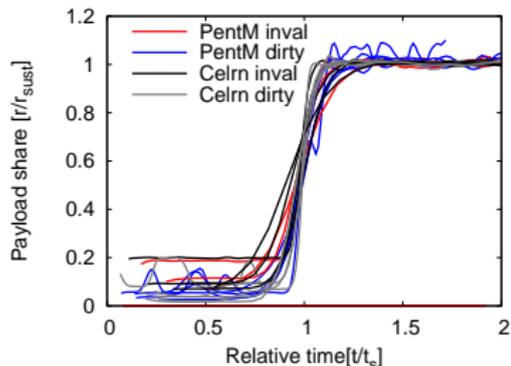
- Messergebnisse (i386):

Machine	WSS	Cache	f_0	t_{sust}
Cel.@2.5G	16k	dirty	0.11	22 μ s
Cel.@2.5G	32k	dirty	0.11	43 μ s
Cel.@2.5G	64k	dirty	0.09	85 μ s
Cel.@2.5G	16k	invd	0.17	12 μ s
Cel.@2.5G	32k	invd	0.18	20 μ s
Cel.@2.5G	64k	invd	0.20	35 μ s
PentM@1.5G	16k	dirty	0.09	25 μ s
PentM@1.5G	32k	dirty	0.08	38 μ s
PentM@1.5G	64k	dirty	0.12	91 μ s
PentM@1.5G	16k	invd	0.17	16 μ s
PentM@1.5G	32k	invd	0.12	29 μ s
PentM@1.5G	64k	invd	0.21	55 μ s

Testprozedur: Aufsteigende Cache-Zeilen
invd = Cache invalidated, dirty = cache flood-filled

- t_{sust} Wächst mit „working space“
- f_0 Zwischen (hier) 8% und 21%

- Normalisierte Werte:



⇒ Entspricht in guter Näherung dem erwarteten Verhalten

- Messergebnisse (PowerPC^a):
 - Single-level Cache, 16 KB

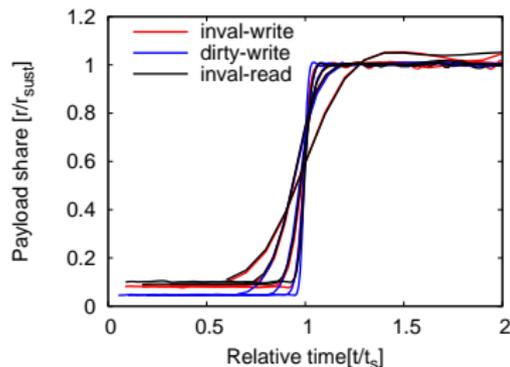
Testcase	WSS	Cache	f_0	t_{sust}
consec_wr	2k	dirty	0.05	17 μ s
consec_wr	4k	dirty	0.05	31 μ s
consec_wr	8k	dirty	0.05	59 μ s
consec_wr	16k	dirty	0.05	116 μ s
consec_wr	2k	invd	0.13	10 μ s
consec_wr	4k	invd	0.09	19 μ s
consec_wr	8k	invd	0.08	35 μ s
consec_wr	16k	invd	0.08	69 μ s
consec_rd	2k	invd	0.13	10 μ s
consec_rd	4k	invd	0.10	19 μ s
consec_rd	8k	invd	0.09	35 μ s
consec_rd	16k	invd	0.10	68 μ s

invd = Cache invalidated, dirty = cache flood-filled

- $t_s \sim WSS$
- f_0 zwischen 5% und 13%
unabhängig von WSS

^aMPC 5200 @ 400MHz

- Normalisierte Werte:

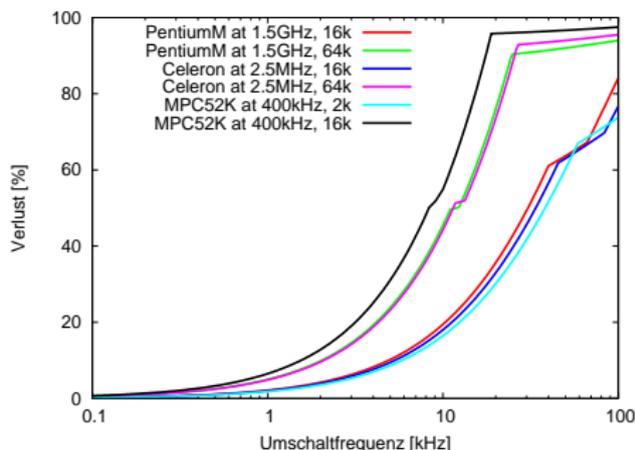


⇒ Entspricht Modellverhalten

Simulation eines Schedulers, parametrisiert mit ermittelten Werten

→ Worst-case-Umschaltverluste in Abhängigkeit von der Umschaltfrequenz

- Unterhalb ca. 1kHz: < 1%
→ vernachlässigbar
- Darüber: Drastischer Anstieg
(bis über 90%)



Faustregel: Zeitscheibendauern unter einer Millisekunde führen zu hohen Verlusten

- Verluste können Ausführungszeiten zugeschlagen werden

1 Einleitung

- Motivation
- Herausforderungen

2 Planung unter Virtualisierungsumgebungen

- Hierarchischer Scheduler
- Planbarkeit

3 Kosten der Virtualisierung

- Planungsbedingte Verluste
- Umschaltkosten

4 Zusammenfassung

- Stellvertreterprozess-Modell zur Zusammenfassung von nach RMS priorisierten Prozessmengen
- Einfach anzuwendendes Kriterium zum Nachweis der Planbarkeit solcher Prozessmengen in virtuellen Maschinen
- Planungsbedingte Verluste...
 - + ... sind beschränkt
 - ... können bis zu 50% betragen
 - Selbst dann: ökonomisch sinnvoll
- Modell zur Approximation Cache-bedingter Umschaltkosten
 - + Bei Frequenzen unter 100 Hz → vernachlässigbar
 - Bei Frequenzen über ca 10kHz → erheblich
 - Kompromiss unter Berücksichtigung der Anwendungsanforderungen

Danke für Ihre Aufmerksamkeit!

... Fragen?