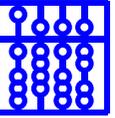

Laufzeit-Werkzeuge für parallele und verteilte Systeme

Roland Wismüller

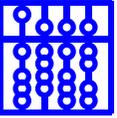
LRR-TUM
Technische Universität München



Werkzeuge zur Laufzeitanalyse

- ➔ *Beobachten, Analysieren* und *Steuern* des Ablaufs eines verteilten Hard- und Softwaresystems (Zielsystem)
- ➔ Hier betrachtet: *on-line* Werkzeuge
- ➔ Beispiele:
 - ➔ Quellsprach-Debugging
 - ➔ Visualisierung des Programmablaufs
 - ➔ Deterministische Ablaufsteuerung
 - ➔ Leistungsanalyse
 - ➔ Sicherungspunkterstellung und Prozeßmigration
 - ➔ Lastausgleich

1. Einführung
2. **Werkzeug-Entwicklung am LRR-TUM**
3. **Online-Monitoring**
4. **Interoperable Werkzeuge**
5. **Zusammenfassung und Ausblick**



Verteilter Debugger DETOP

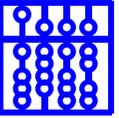
➔ Realisiert für MMK, PARIX, PVM, MPI, HPF (SPiDER)

Fragestellungen

- ➔ Effiziente Mechanismen zum Monitoring von Threads
- ➔ Skalierbarkeit
- ➔ Compiler-Optimierungen
- ➔ Mehrstufige Übersetzung / Parallelisierung

The screenshot shows a debugger window titled "DETOP window" with the following components:

- Threads List:** A list of threads on the left, including n0.c3.t1 (delivered), n3.c3.t1, n0.c3.t2, n0.c3.t3, n3.c3.t2, n3.c3.t3, n3.c3.t4, n0.c3.t4, n3.c3.t5, n0.c3.t5, n2.c3.t1, n1.c3.t1 (produced), n1.c3.t10 (produced), n1.c3.t2 (produced), n1.c3.t3 (produced), n1.c3.t4 (produced), n1.c3.t5 (produced), n1.c3.t6 (produced), n1.c3.t7 (produced), n1.c3.t8 (produced), n1.c3.t9 (produced), n2.c3.t6, n2.c3.t7, and n2.c3.t8.
- File: tst1.c:** The main source code window showing lines 41-57. A green arrow points to line 55, indicating the current execution point. Two red "STOP" icons mark breakpoints at lines 45 and 57.
- State:** A panel on the right showing the state of thread n1.c3.t1 as "STOPPED". It lists the call stack: produce at tst1.c:52, _PX_usrfunc_harness (no line info), and ParixMain (no line info).
- Output:** A panel on the right showing the output of a print statement: "Value of `tst1.c:produce` buffer is: [n1.c3.t1]". The buffer content is displayed as a struct with fields: Header (Size: 6542804, DestProcId: 8160184, ReqId: 1395076144, MsgType: 808858937, SourceProcId: 6507976, Code: 658843436, Timeout: 33) and Body (array[1024]).
- Command line:** A field at the bottom containing "n0.*(SendProc,*)".



Compiler-Optimierungen

Quellprogramm

10 v = ...

11 w = ...

12 v = ...

Zielprogramm

1234: store ...,r0

...

1236: store ...,r0

...

1239: store ...,r1

Compiler-Optimierungen

Quellprogramm

```

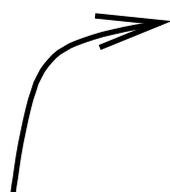
10 v = ...
11 w = ...
12 v = ...

```

```

break 12
print v, w

```



Zielprogramm

```

1234: store ...,r0
...
1236: store ...,r0
...
1239: store ...,r1

```

Compiler-Optimierungen

Quellprogramm

Zielprogramm

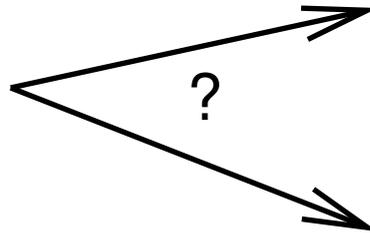
```

10 v = ...
11 w = ...
12 v = ...
  
```

```

1234: store ...,r0
...
1236: store ...,r0
...
1239: store ...,r1
  
```

break 12
print v, w



1. Abbildung des Haltepunkts?

Compiler-Optimierungen

Quellprogramm

```

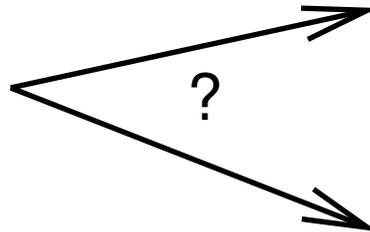
10 v = ...
11 w = ...
12 v = ...
  
```

break 12
print v, w

Zielprogramm

```

1234: store ...,r0
...
1236: store ...,r0
...
1239: store ...,r1
  
```



1. Abbildung des Haltepunkts?
2. Verfügbarkeit von Variablenwerten?

Compiler-Optimierungen

Quellprogramm

```

10 v = ...
11 w = ...
12 v = ...

```

break 12
print v, w

Zielprogramm

```

1234: store ...,r0
...
1236: store ...,r0
...
1239: store ...,r1

```

1. Abbildung des Haltepunkts?
2. Verfügbarkeit von Variablenwerten?
 - r0 enthält Wert von v
 - Wert von w nicht verfügbar

Compiler-Optimierungen

Quellprogramm

```

10 v = ...
11 w = ...
12 v = ...

```

```

break 12
print v, w

```

Zielprogramm

```

1234: store ...,r0
...
1236: store ...,r0
...
1239: store ...,r1

```

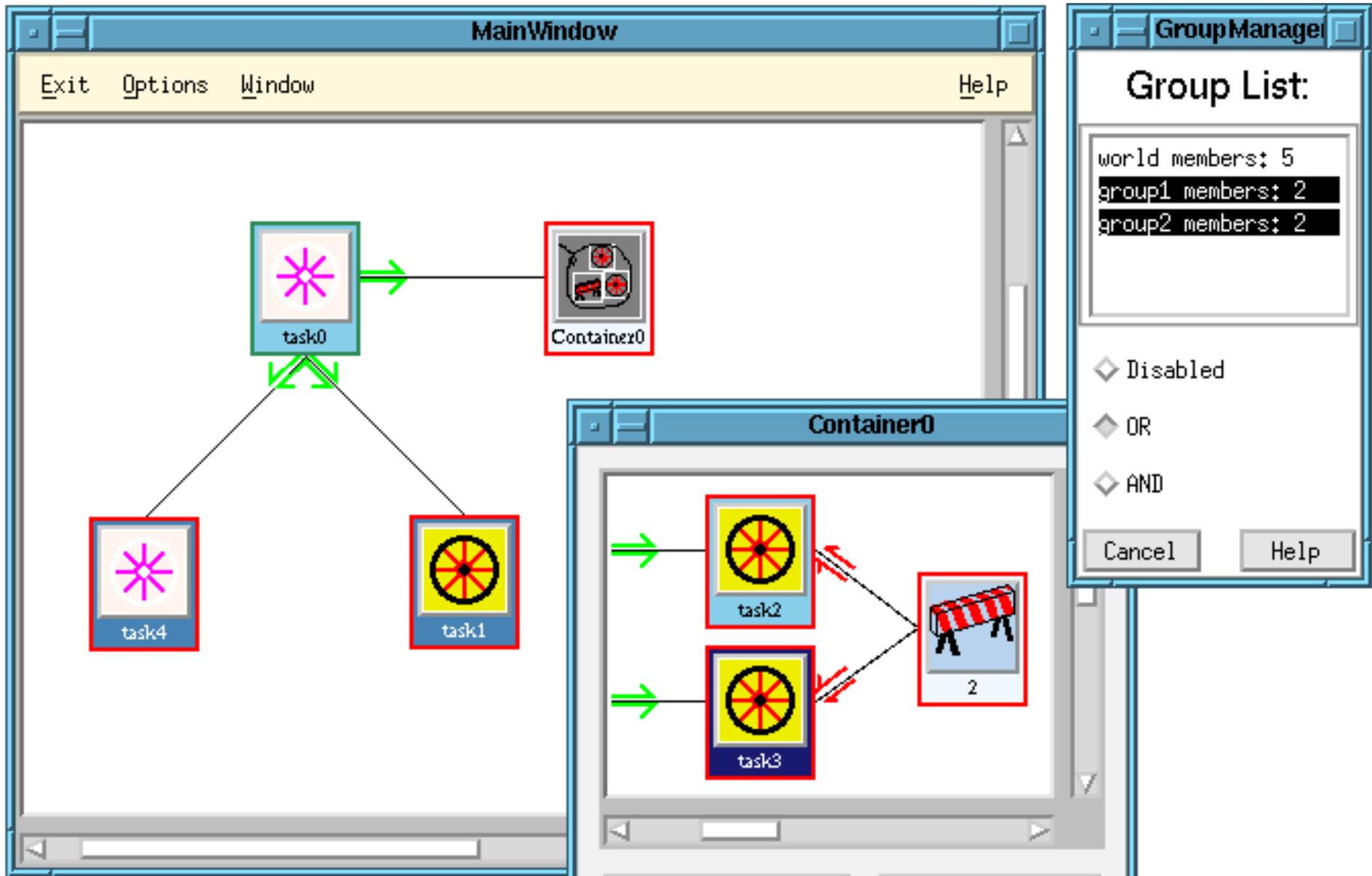
1. Abbildung des Haltepunkts?
2. Verfügbarkeit von Variablenwerten?
 - r1 enthält Wert von w
 - Wert von v nicht verfügbar

Visualisierer VISTOP

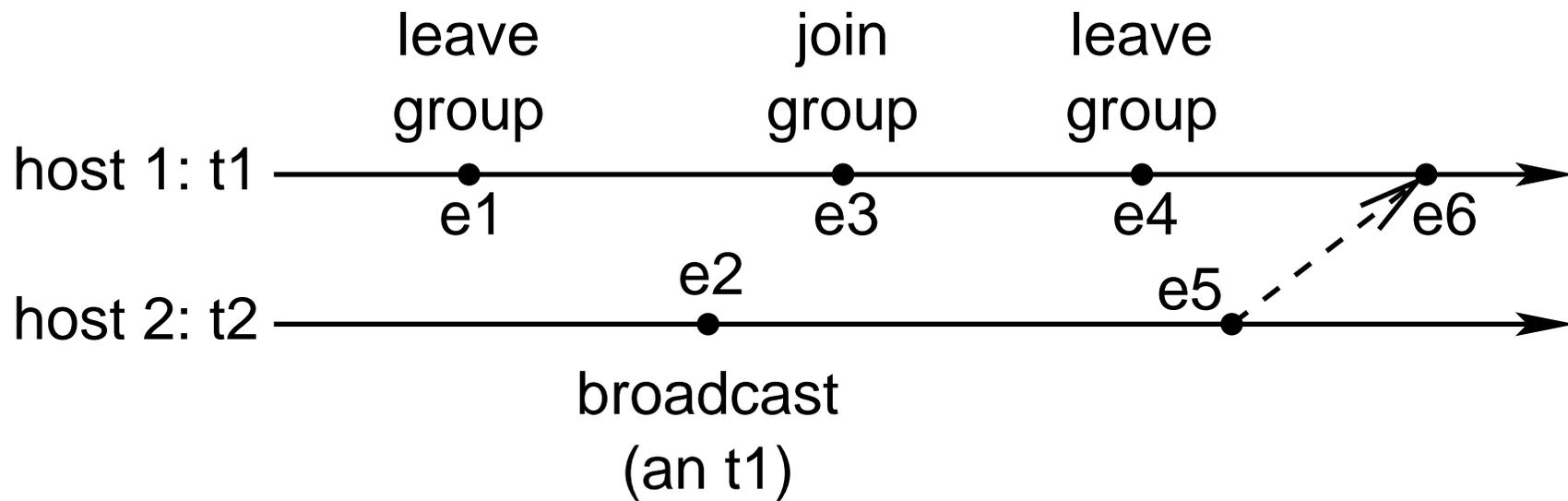
- ➔ zeigt (mögliche) Folgen globaler Zustände
 - ➔ Kommunikation, Synchronisation, Prozeßgruppen
- ➔ Realisiert für MMK und PVM

Fragestellungen

- ➔ Visualisierungsmetaphern
- ➔ Datenerfassungstechniken
- ➔ Berechnung konsistenter globaler Zustände

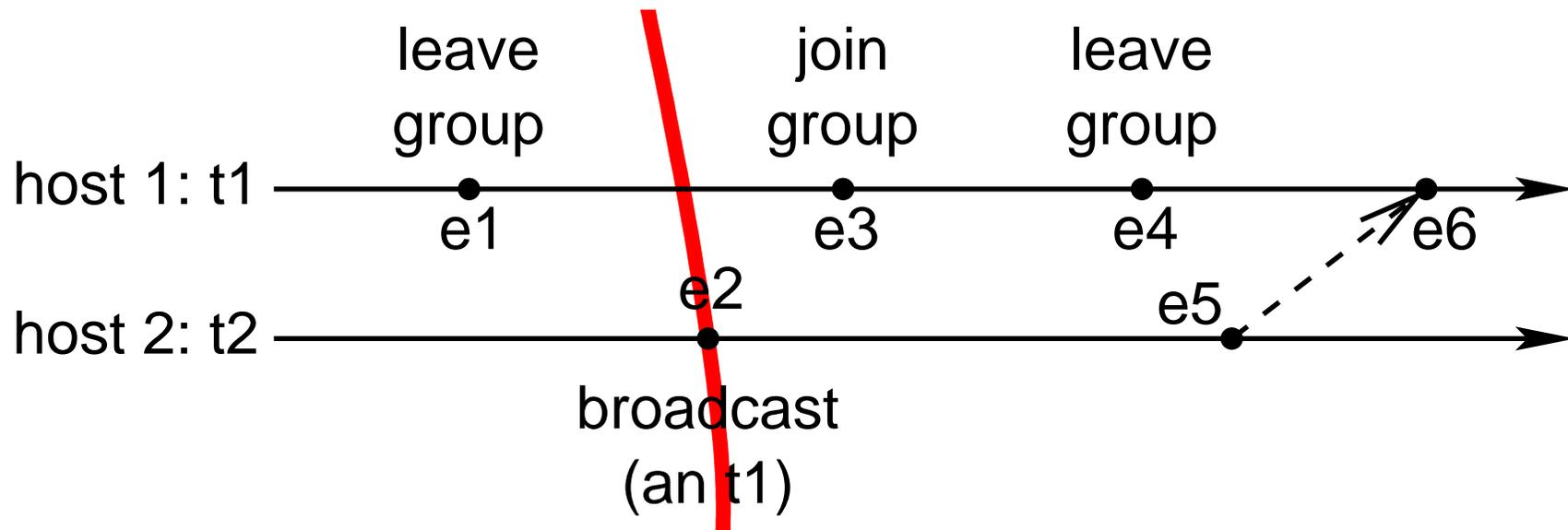


Berechnung konsistenter globaler Zustände



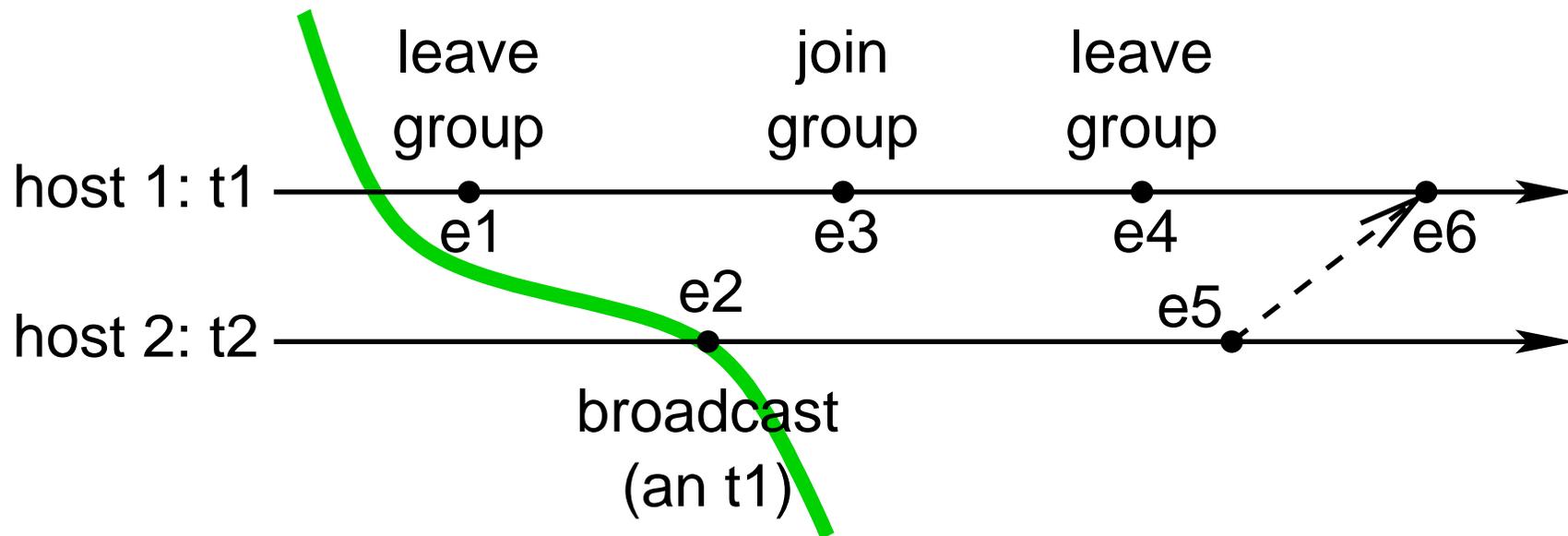
- ➔ Reihenfolge $e_1, e_2, e_3, e_4, e_5, e_6$ offenbar inkonsistent
- ➔ Reihenfolgen $e_2, e_1, e_3, e_4, \dots$ und $e_1, e_3, e_2, e_4, \dots$ konsistent
- ➔ Keine Kausalitätsbeziehung zwischen e_2, e_1 bzw. e_2, e_3

Berechnung konsistenter globaler Zustände



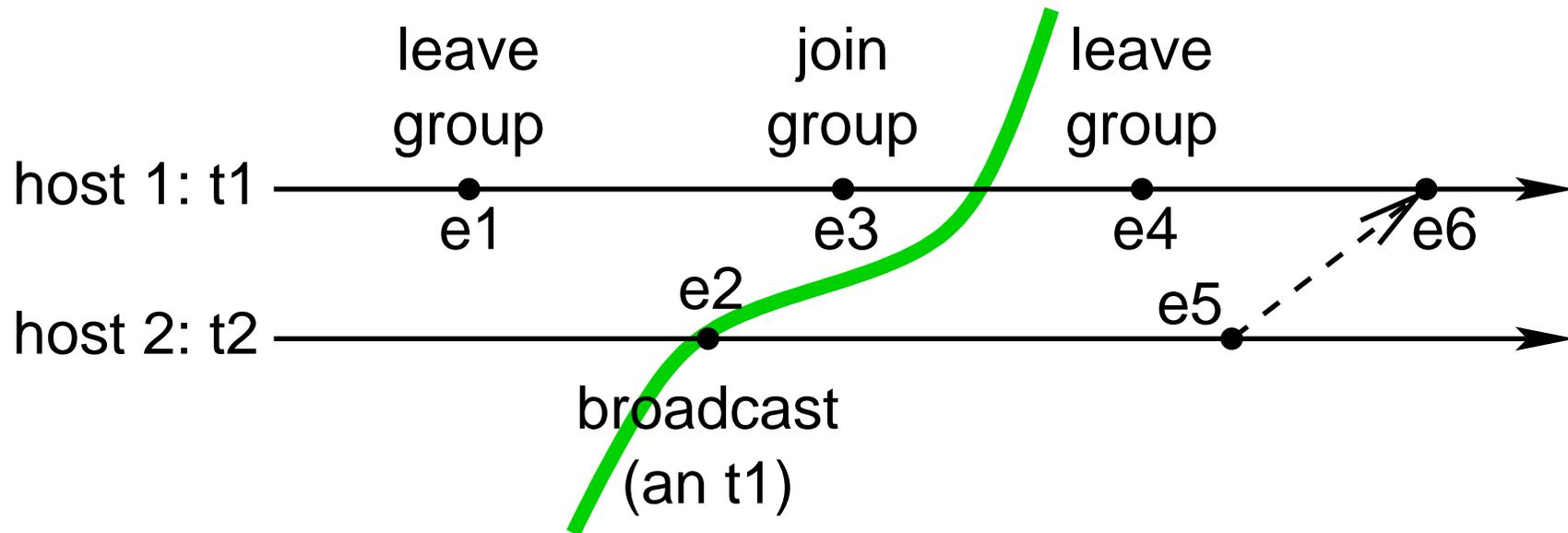
- ➔ Reihenfolge $e_1, e_2, e_3, e_4, e_5, e_6$ offenbar inkonsistent
- ➔ Reihenfolgen $e_2, e_1, e_3, e_4, \dots$ und $e_1, e_3, e_2, e_4, \dots$ konsistent
- ➔ Keine Kausalitätsbeziehung zwischen e_2, e_1 bzw. e_2, e_3

Berechnung konsistenter globaler Zustände



- ➔ Reihenfolge $e1, e2, e3, e4, e5, e6$ offenbar inkonsistent
- ➔ Reihenfolgen $e2, e1, e3, e4, \dots$ und $e1, e3, e2, e4, \dots$ konsistent
- ➔ Keine Kausalitätsbeziehung zwischen $e2, e1$ bzw. $e2, e3$

Berechnung konsistenter globaler Zustände



- ➔ Reihenfolge $e1, e2, e3, e4, e5, e6$ offenbar inkonsistent
- ➔ Reihenfolgen $e2, e1, e3, e4, \dots$ und $e1, e3, e2, e4, \dots$ konsistent
- ➔ Keine Kausalitätsbeziehung zwischen $e2, e1$ bzw. $e2, e3$

Ziel: Testunterstützung

- ➔ Nicht nur *Replay*
- ➔ Erzwingen möglicher Abläufe durch Zugriffsmuster für gemeinsame Objekte

Ablaufsteuerungs-Werkzeug codex

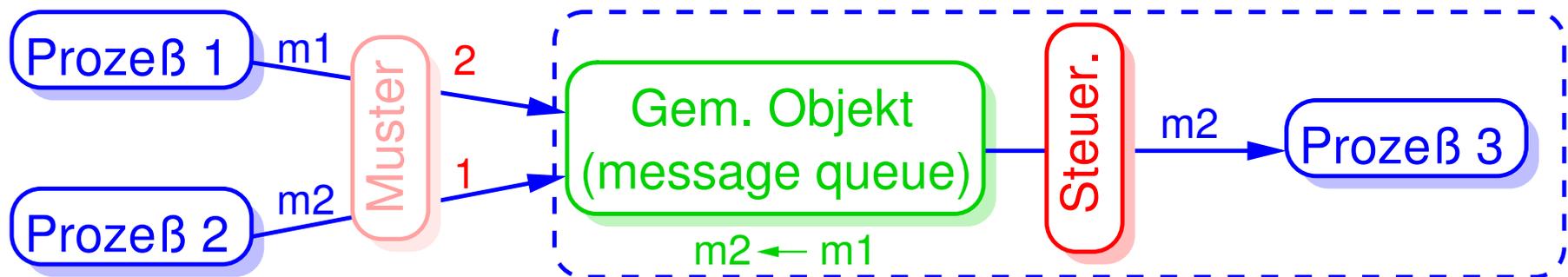
- ➔ Realisiert für PVM

Fragestellungen

- ➔ Spezifikation von Zugriffsmustern
- ➔ Effizientes Erzwingen der Zugriffsmuster

Effizientes Erzwingen der Zugriffsmuster

- ➔ Für PVM: Spezielle Semantik der gemeinsamen Objekte (*message queues*)
- ➔ Konzeptuell: Muster steuert Reihenfolge der Schreibzugriffe (Sendeaufrufe)
- ➔ Realisierung: Umordnung der Nachrichten erst beim Empfang



Leistungsanalyse-Werkzeuge PATOP und TATOO

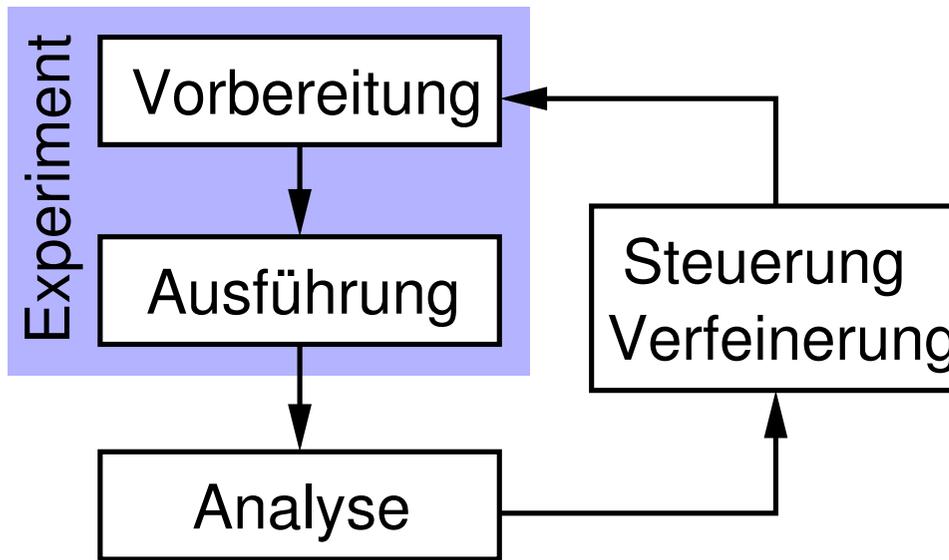
- ➔ Online- und Offline-Analyse
- ➔ Realisiert für MMK, PARIX, PVM, MPI, HPF, ...

Fragestellungen

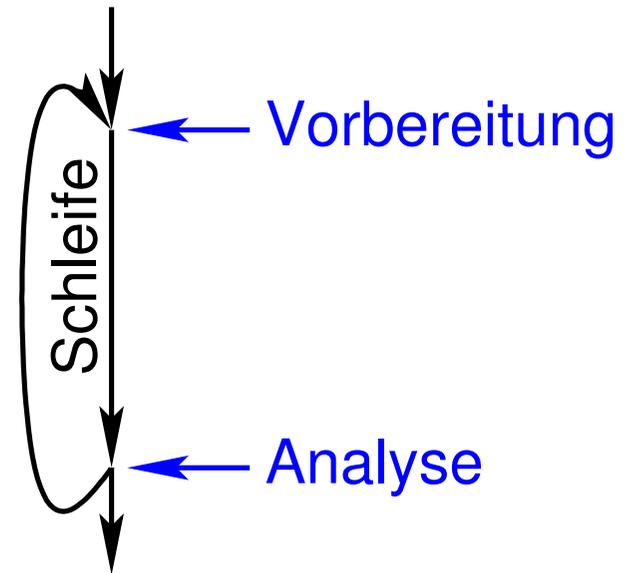
- ➔ Online-Analyse mit verteilter Auswertung
- ➔ Leistungsanalyse in nicht-exklusiven Systemen
- ➔ Automatische Online-Leistungsanalyse

Automatische Online-Leistungsanalyse

Performance-Experiment



Online-Experiment



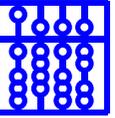
- ➔ Verfeinerung von Hypothesen nur für wiederholte Phasen möglich
- ➔ Werkzeug benötigt Strukturinformation (Programmierer, Compiler)

Werkzeuge CoCheck, LoMan, LMC

- ➔ CoCheck: Sicherungspunkte und Migration für PVM
- ➔ LoMan: Lastverwalter auf Basis von CoCheck
- ➔ LMC: Lastverwaltung für CORBA

Fragestellungen

- ➔ Protokolle (konsistente Sicherungspunkte, Migration)
- ➔ Lastmessung und -bewertung
- ➔ Kombination von Initialplatzierung, Migration und Replikation

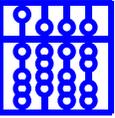


Motivation

- ➔ Jedes Online-Werkzeug benötigt verteilte Infrastruktur zur
 - ➔ Beobachtung
 - ➔ Manipulation
 - von Programmabläufen

- ➔ Bisher:
 - ➔ Werkzeuge besitzen spezielle, proprietäre Monitorsysteme
 - ➔ Keine Standard-Schnittstellen

- ➔ Probleme:
 - ➔ Entwicklungs- und Portierungsaufwand
 - ➔ Interoperabilität



Anforderungen an eine Standard-Schnittstelle

- ➔ Unterstützung für Online-Monitoring
- ➔ Plattformunabhängigkeit
- ➔ Nutzbarkeit für verschiedenste Werkzeuge
- ➔ Unterstützung interoperabler Werkzeuge
- ➔ Effiziente, skalierbare Implementierung

Umsetzung:

- ➔ OMIS (On-line Monitoring Interface Specification)

Konzepte von OMIS

- ➔ Objektbasiertes Modell des Zielsystems
 - ➔ Objekte mit abstrakten Identifikatoren
 - ➔ Dienste auf diesen Objekten
 - ➔ Werkzeuge definieren ihre Sicht des Zielsystems
- ➔ Ereignis/Aktions-Modell
 - ➔ Anfrage = Ereignisdefinition + Aktionsliste
 - ➔ Übergabemechanismus für Ereignisinformation
- ➔ Skalierbarkeit: Dienste auf Objektmengen
- ➔ Ortstransparenz
- ➔ Erweiterbarkeit

OMIS Beispielanfragen

```
thread_reached_addr([p_14,p_52], 0x124):
    thread_stop([a_])
```

Wenn ein Thread in einem der Prozesse die Adresse 0x124 erreicht, stoppe alle überwachten Threads.

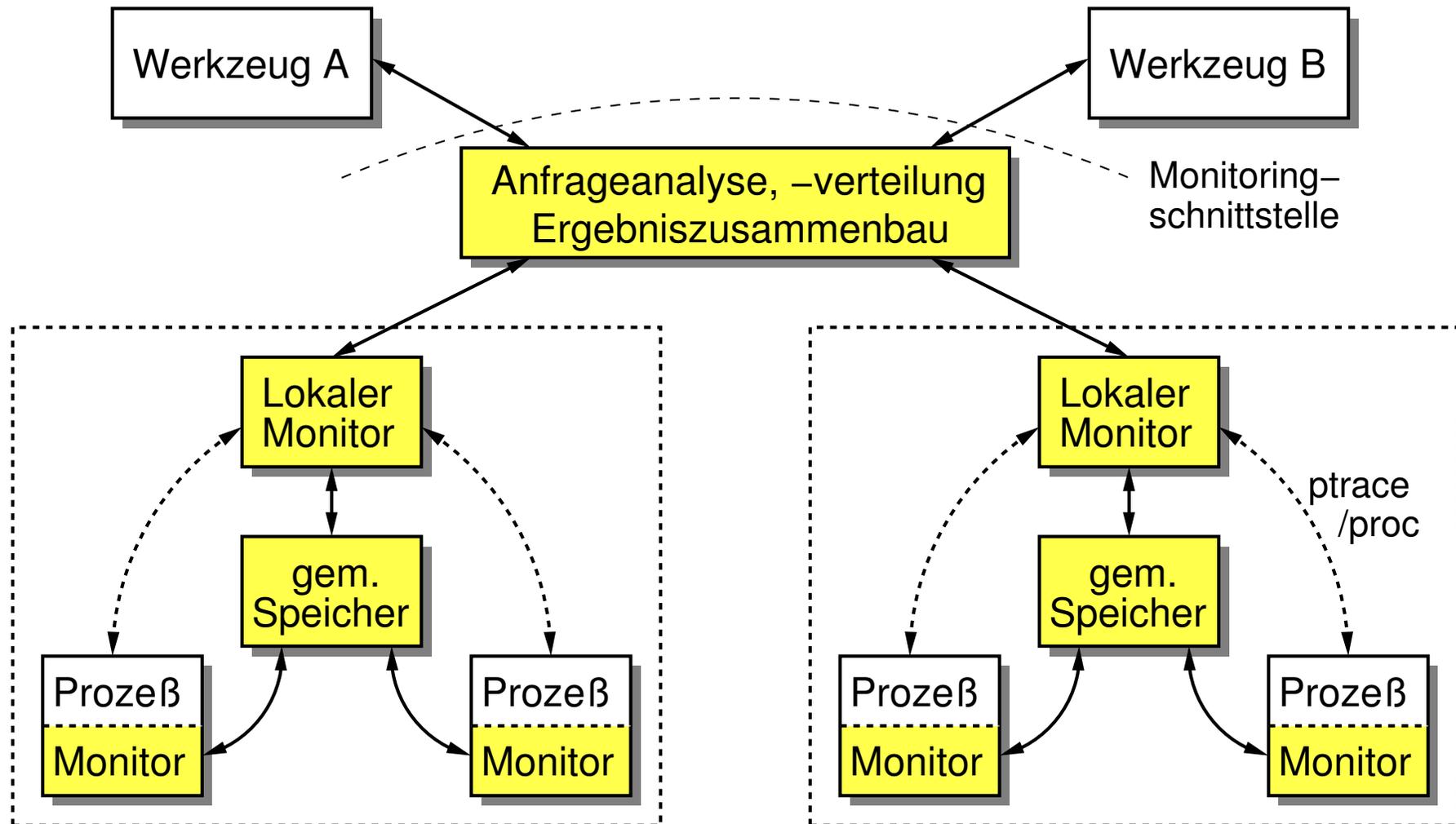
```
thread_has_ended_lib_call([a_], "pvm_recv") :
    pa_counter_global_increment(pa_c_1, 1)
thread_creates_proc([a_]) :
    node_attach([$new_node]) proc_attach([$new_proc])
: pa_counter_global_read([pa_c_1], 1)
```

Anzahl der Aufrufe von `pvm_recv` in allen Threads (Prozessen)

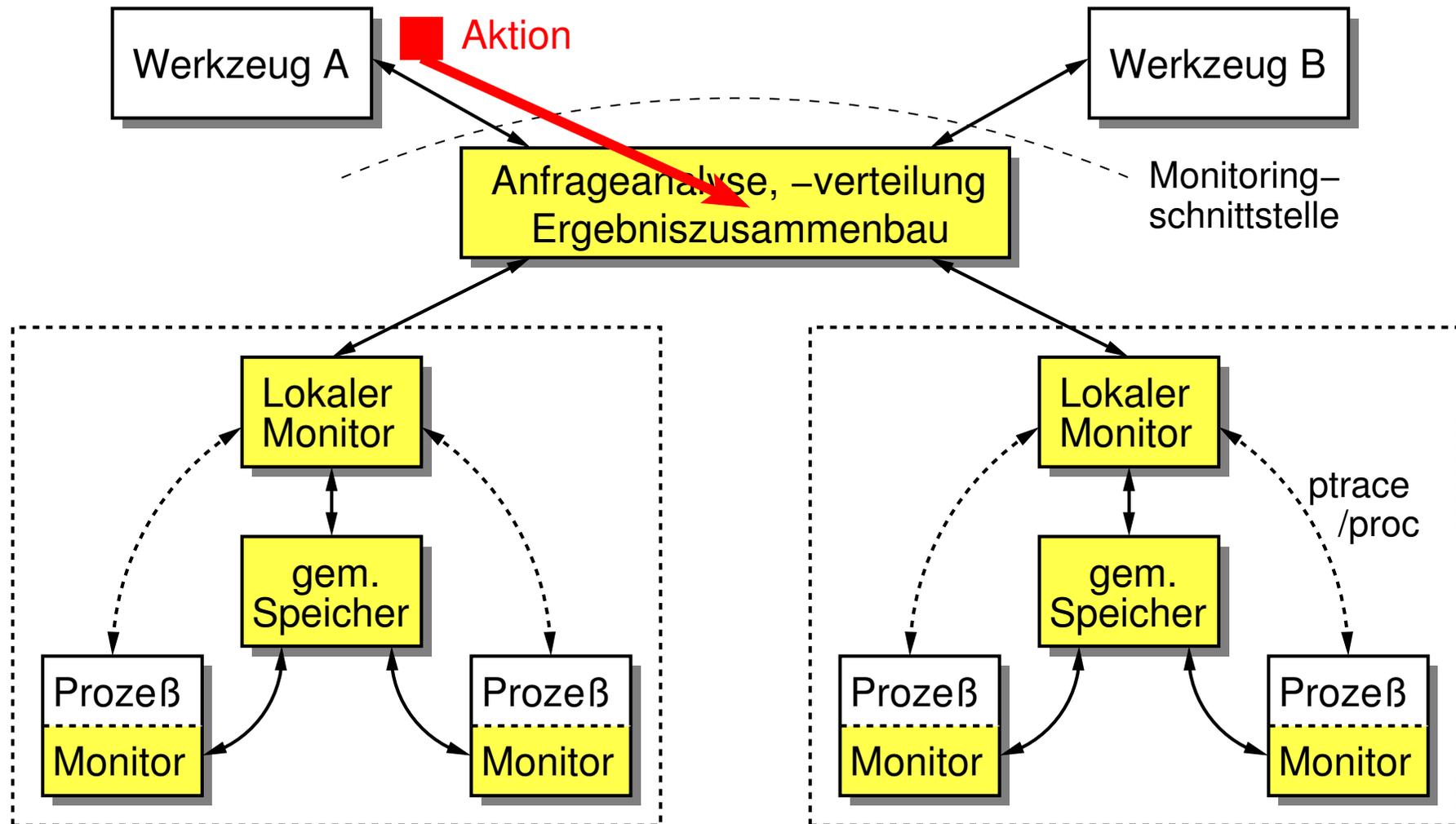
Implementierung von OMIS: OCM

- ➔ Für PVM und MPI auf UNIX-Clustern (Solaris, Irix, Linux, ...)
- ➔ Erste Werkzeuge: DETOP, VISTOP, CoCheck, PATOP
- ➔ Portierungen auf andere Plattformen
 - ➔ DSM-Systeme (SMiLE-Cluster am LRR-TUM, HW-Monitor)
 - ➔ SMP-Systeme (Thread-Monitoring)
 - ➔ Java / RMI
 - ➔ Grid-Umgebungen

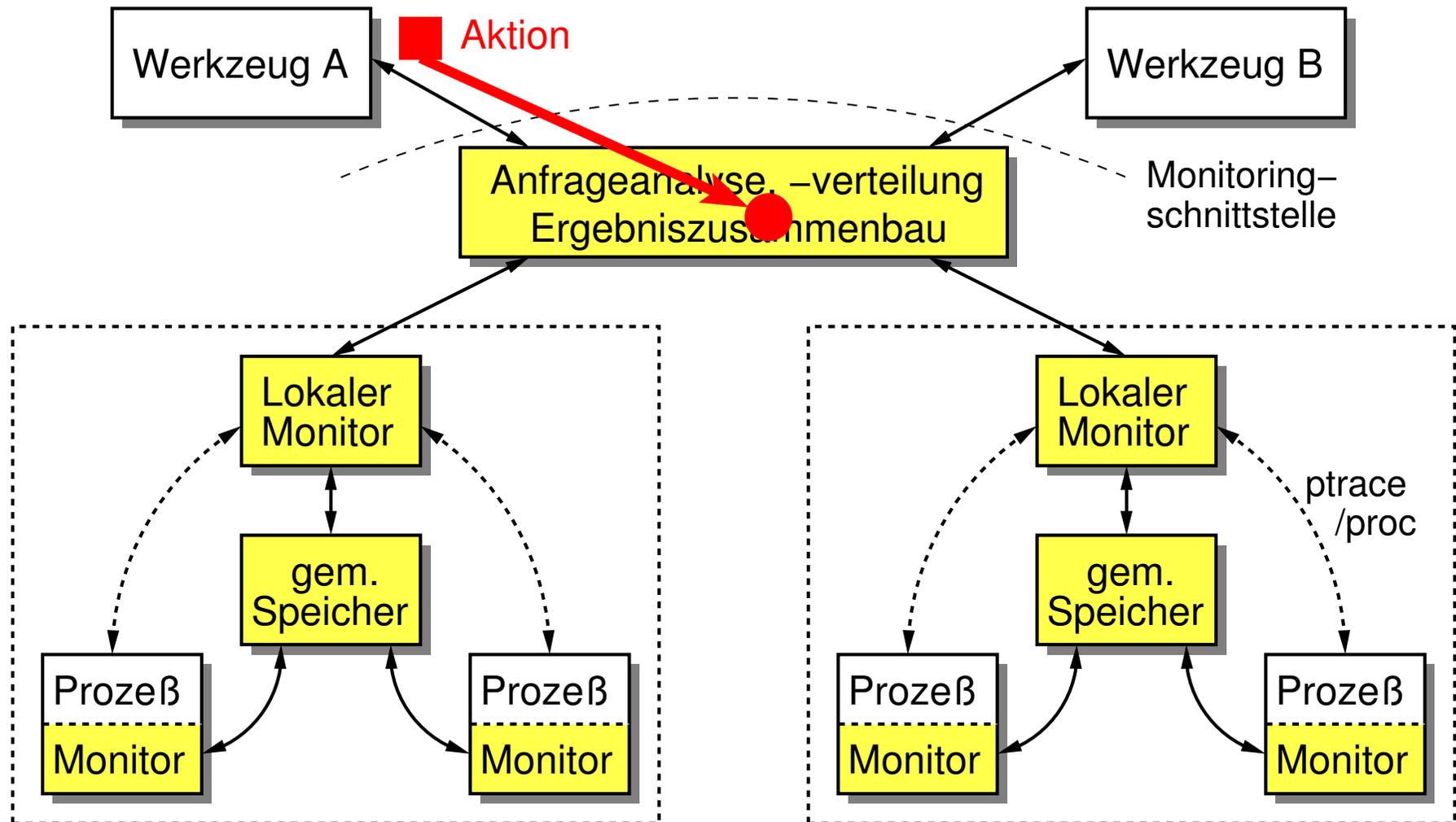
Struktur des OCM



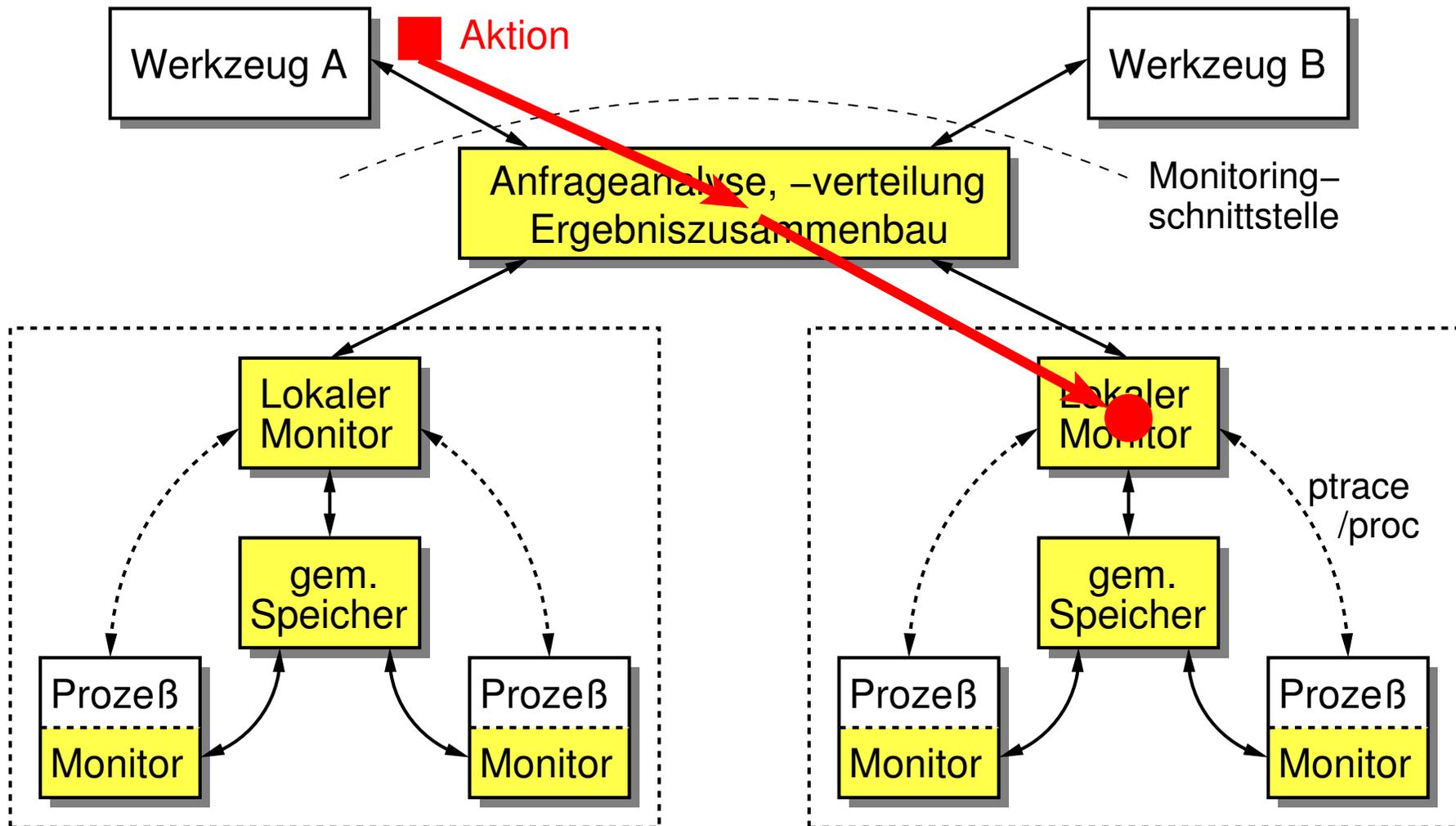
Anfragebearbeitung im OCM



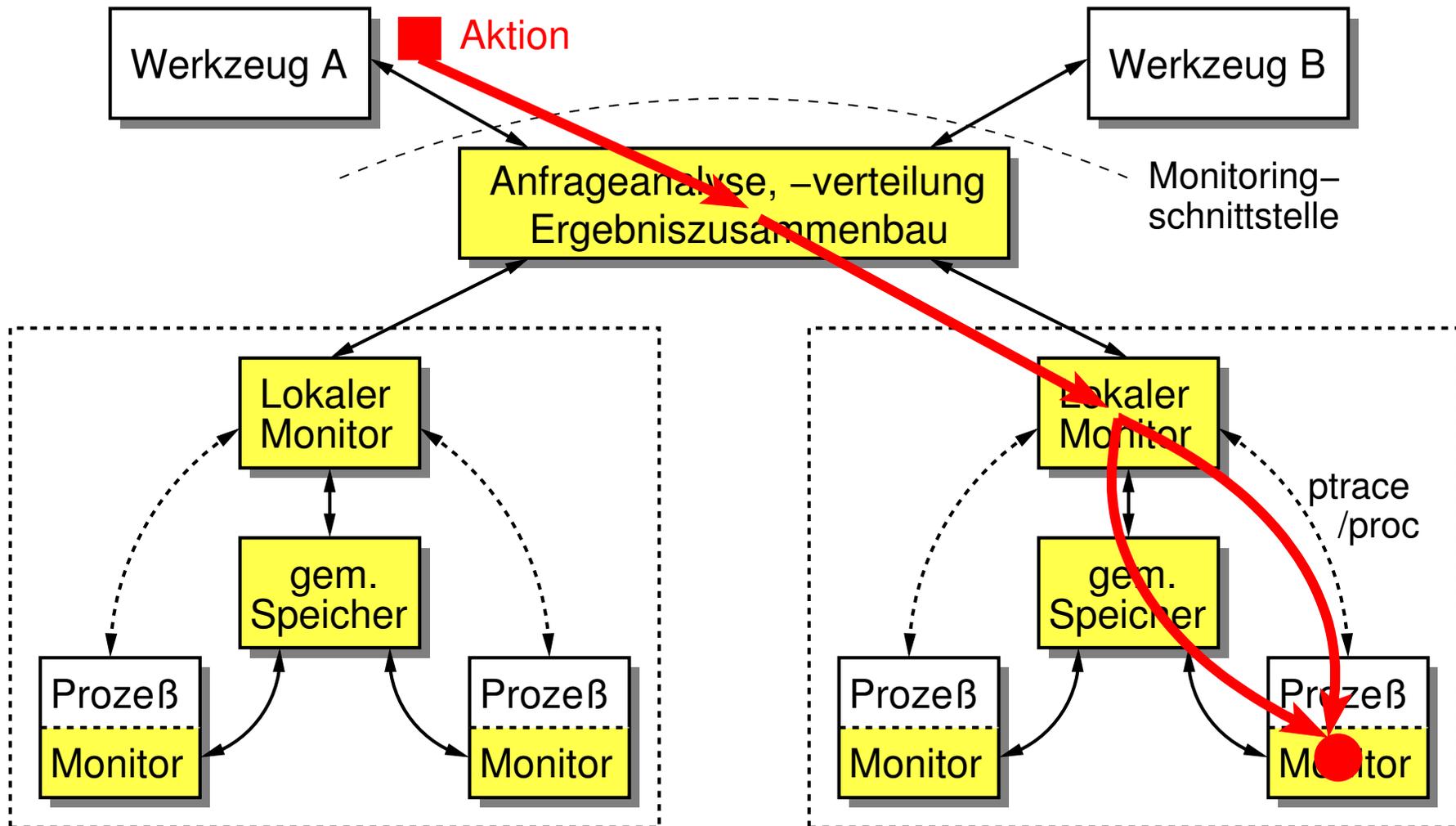
Anfragebearbeitung im OCM



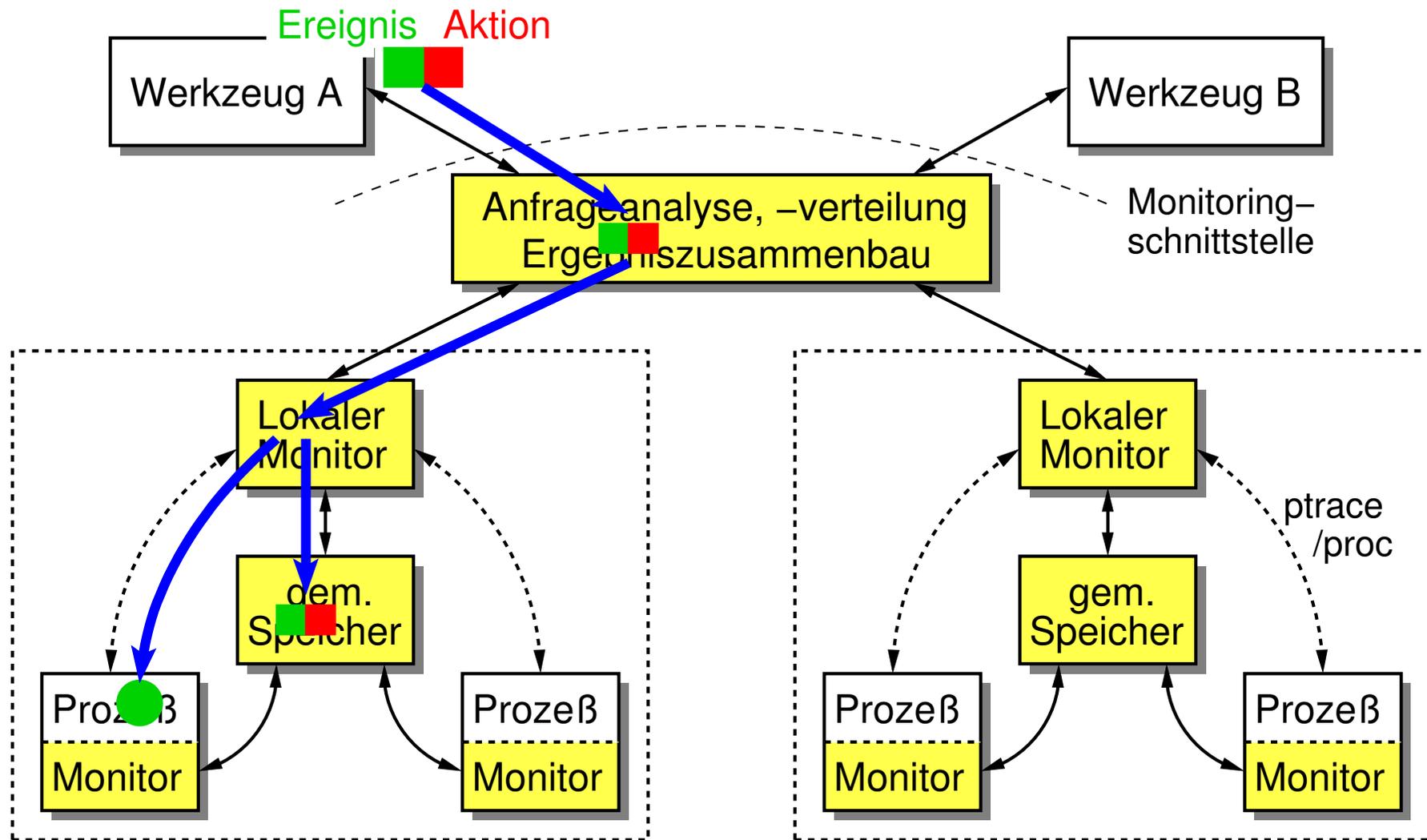
Anfragebearbeitung im OCM



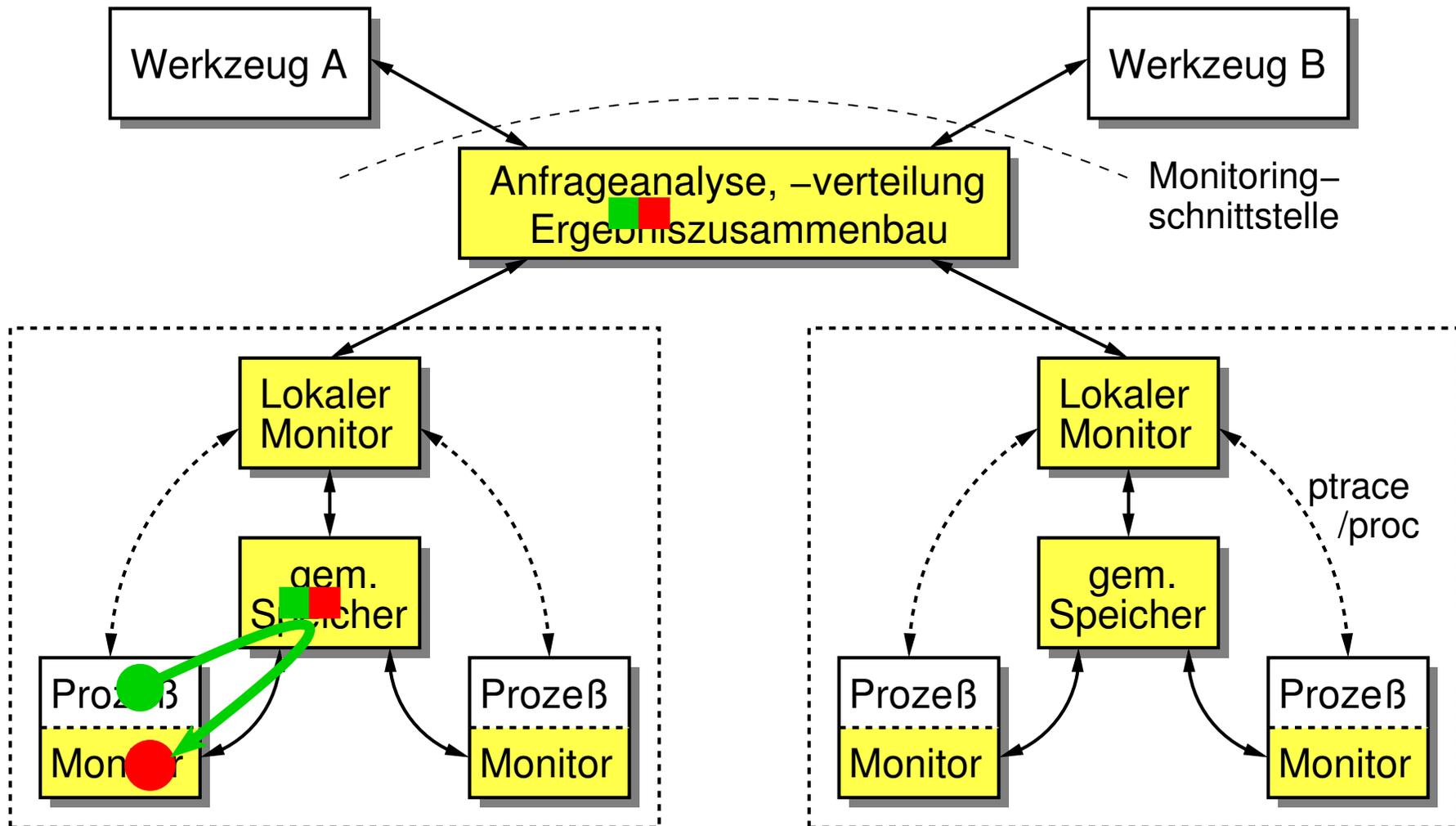
Anfragebearbeitung im OCM



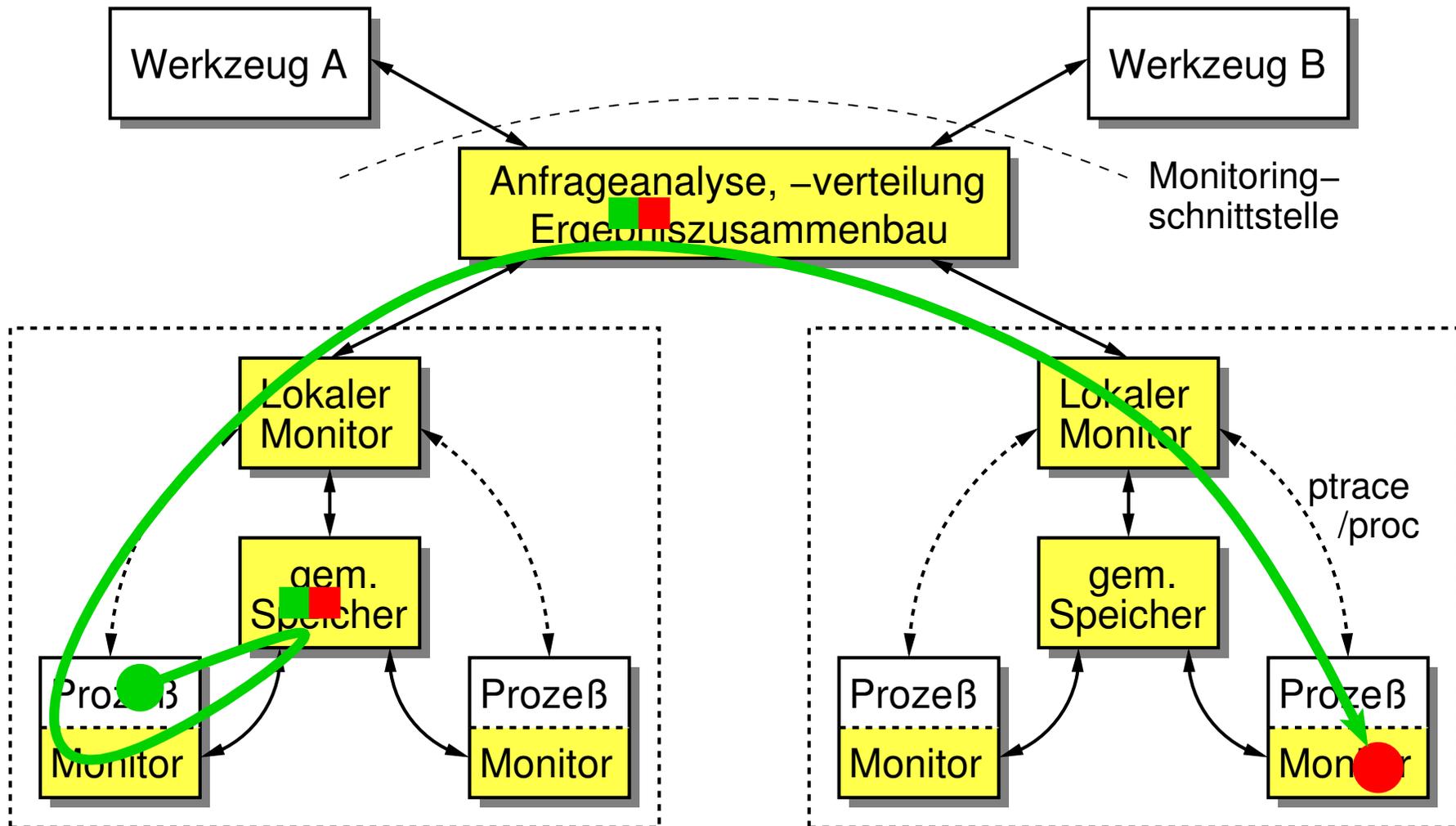
Anfragebearbeitung im OCM



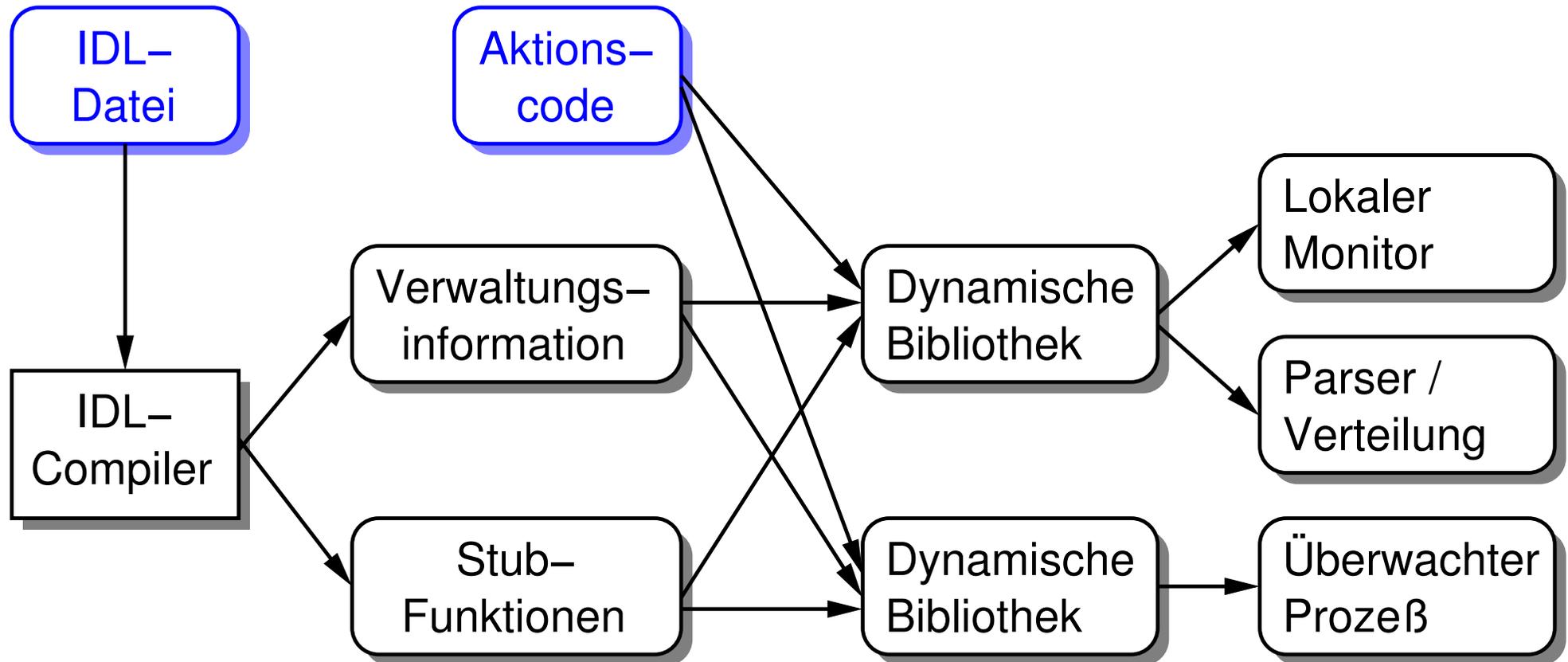
Anfragebearbeitung im OCM

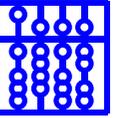


Anfragebearbeitung im OCM



Erweiterung des OCM um neue Aktionen



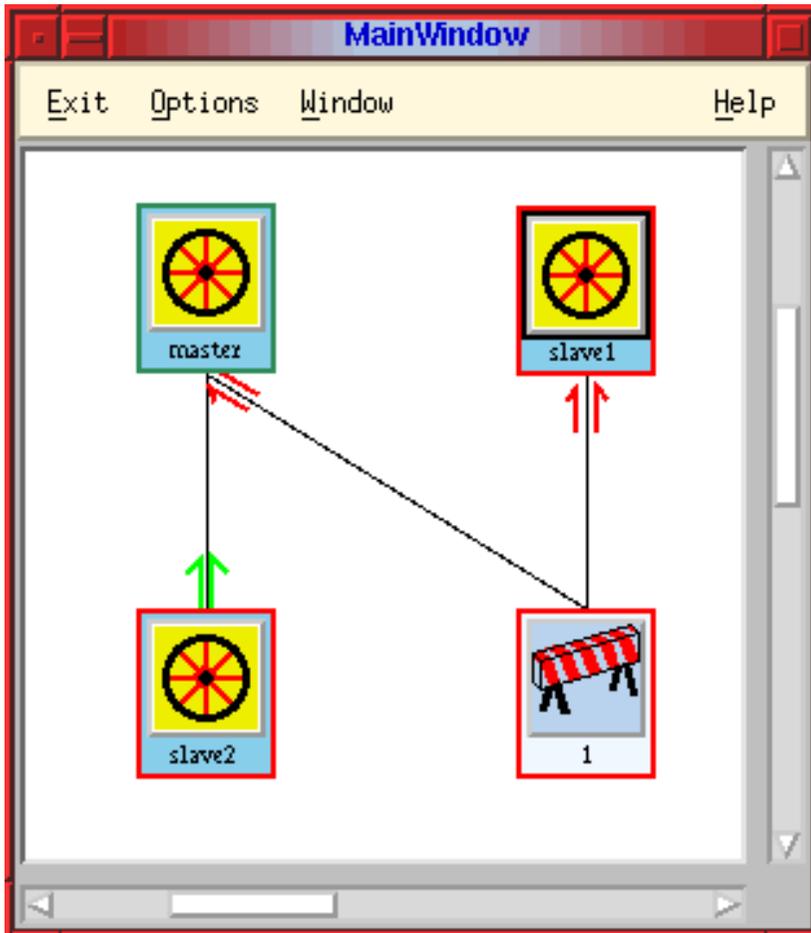


Motivation: Kombination von Werkzeugen

- ➔ z.B. Visualisierer + Debugger
 - ➔ Überblick über Ablaufverhalten
 - ➔ detaillierte Untersuchung

Problem:

- ➔ Daten- und Verhaltensintegration notwendig
- ➔ Keine Abhängigkeit der Werkzeuge
- ➔ Möglichst keine Modifikation der Werkzeuge



Tasks

- p1 (-1, /)
- p2 (2621)
- p3 (2621)

File: appl2.c

```

129
130     printf("*** Slave st
131     pvm_joyngroup("testg
132     pvm_barrier("testgrc
133     for (i=0; i<passes;
134         pvm_recv(-1,1);
135         pvm_unpackf("%d"
136         pvm_packf("%+ %c
137     pvm_send(tid, 2)
138     pvm_barrier("tes
139     }
140     printf("*** Slave er
141     }
142     pvm_exit();
143     exit(0);
144 }
145

```

State

State of task p2:
 State: STOPPED
 Calls:
 main at appl2.c:138
 _start (no line info)

State of task p3:
 State: STOPPED
 Calls:
 main at appl2.c:137
 _start (no line info)

Output

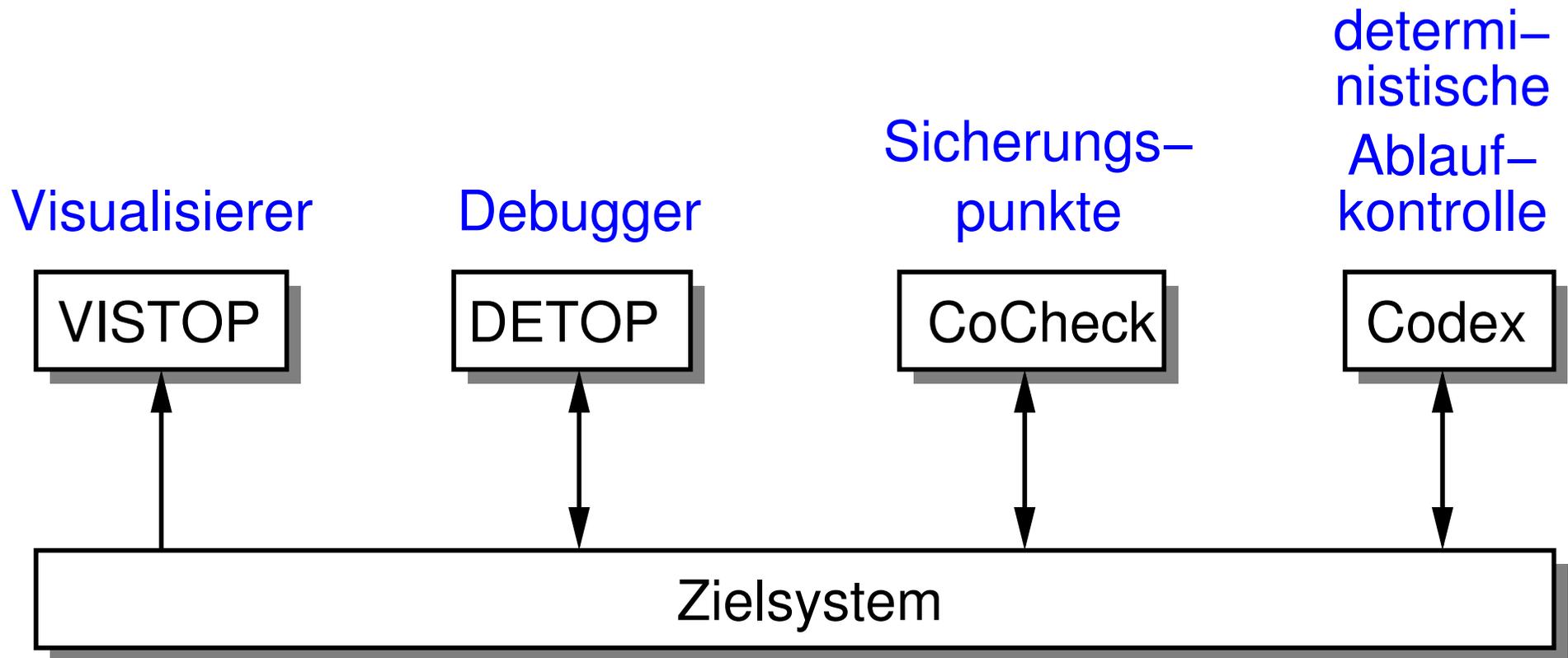
```

> print [p2, p3] `appl2.c:main`
Value of `appl2.c:main` tid is:
[p2, p3]
262149;

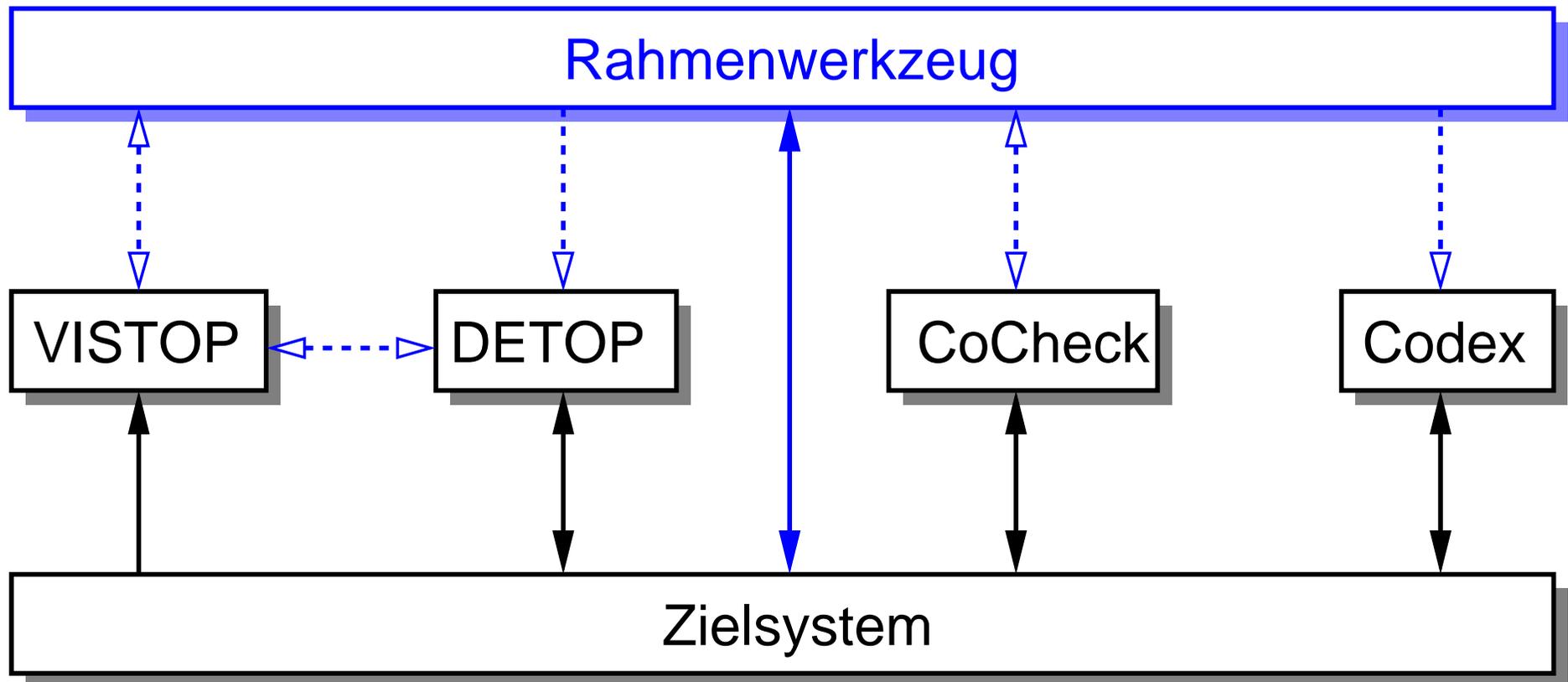
```

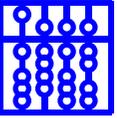
Command line

Aufbau der Werkzeugumgebung



Aufbau der Werkzeugumgebung

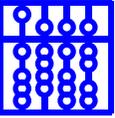




Problem: Konflikte bei Zugriff auf gemeinsames Zielsystem

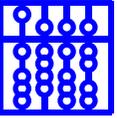
- ➔ Instrumentierung / Betriebssystem-Schnittstellen
 - ➔ Gelöst durch gemeinsames Monitorsystem (OMIS/OCM)

- ➔ Explizite Modifikation des Zielsystemzustands
 - ➔ DETOP hält Zielprozeß an
 - ⇒ VISTOP muß von der Änderung erfahren
 - ➔ CoCheck ändert Task-ID der Prozesse
 - ⇒ andere Werkzeuge dürfen Änderung nicht sehen



Erkennung und Auflösung von Konflikten

- ➔ Lösungsansatz:
 - ➔ Zielsystemobjekt = Menge von Zustandskomponenten mit Schreib-/Leseoperationen
 - ➔ Operationen mit Attribut: Beobachtungsebene
 - ➔ Konflikterkennung und -auflösung anhand der Attribute
- ➔ Konfliktauflösung:
 - ➔ Beobachtung von Schreibzugriffen
 - ➔ Abfangen von Lesezugriffen
- ➔ Zusätzlich berücksichtigt:
 - Überlappung von Zustandskomponenten, Erweiterbarkeit



Ziel: Kooperation zwischen Werkzeugen

- ➔ Z.B. synchronisierte Auswahl, konsistente Darstellung
- ➔ Im weiteren Sinne: Steuerung durch Rahmenprogramm

Allgemeinere Fragestellung:

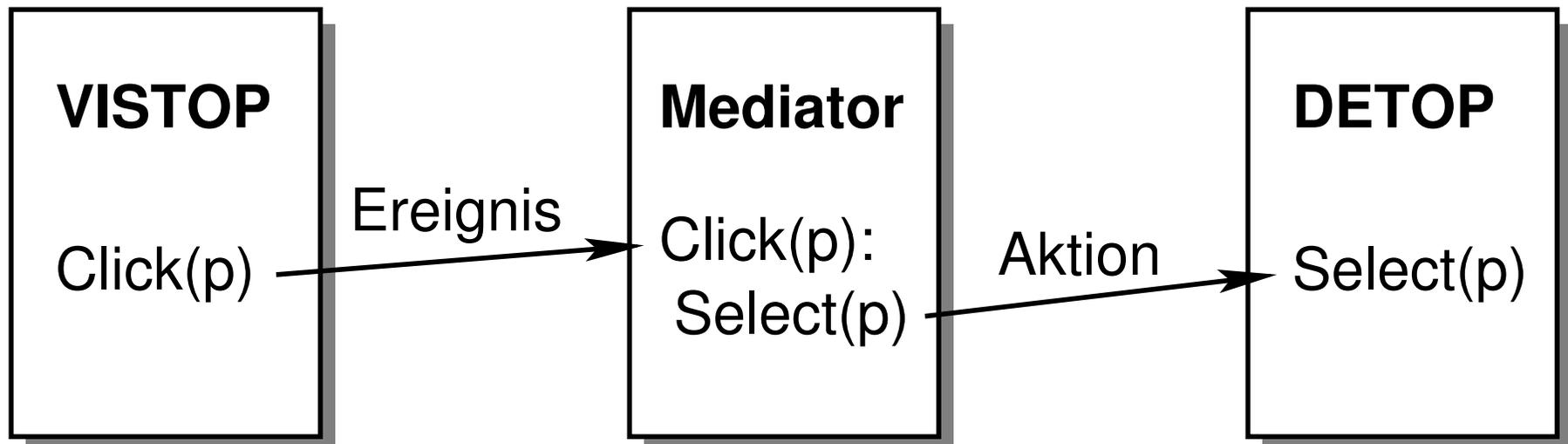
- ➔ Erstellung von Software-Umgebungen aus vorhandenen, aktiven Komponenten (= Programme mit Benutzerinteraktion)

Derzeit bester Ansatz:

- ➔ Mediatoren (bzw. Interaktoren)

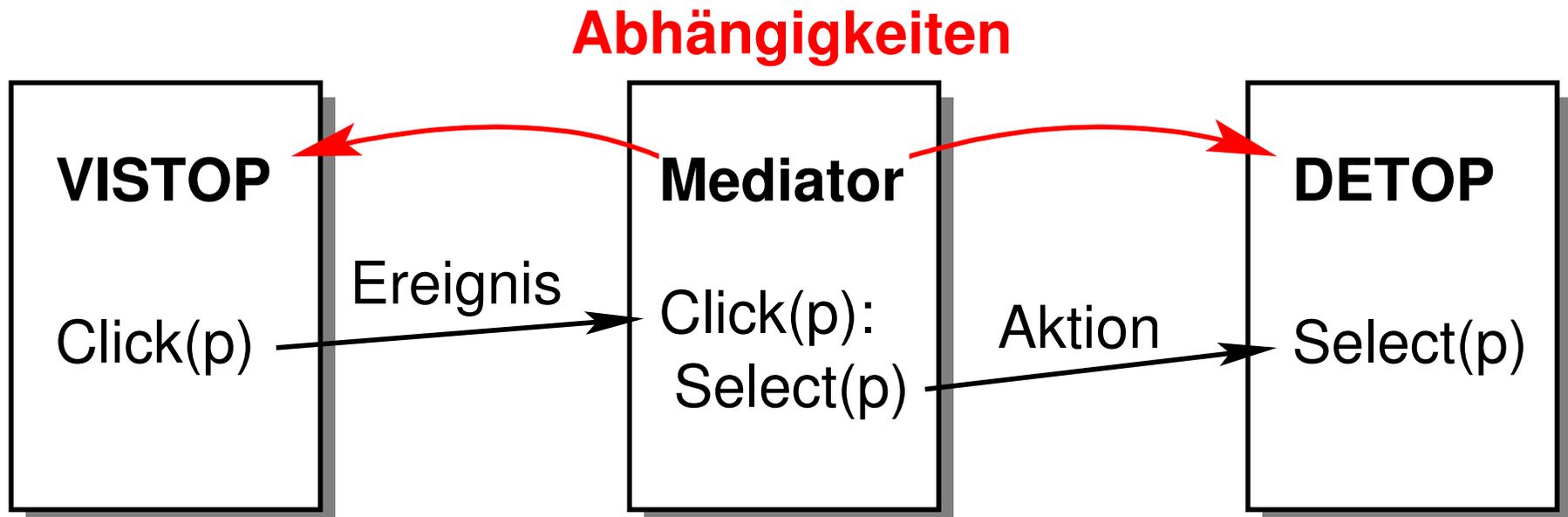
Mediatoren / Interaktoren

- ➔ Mediator ist zusätzliche Komponente
- ➔ Überwacht Ereignisse, führt zugehörige Aktionen aus

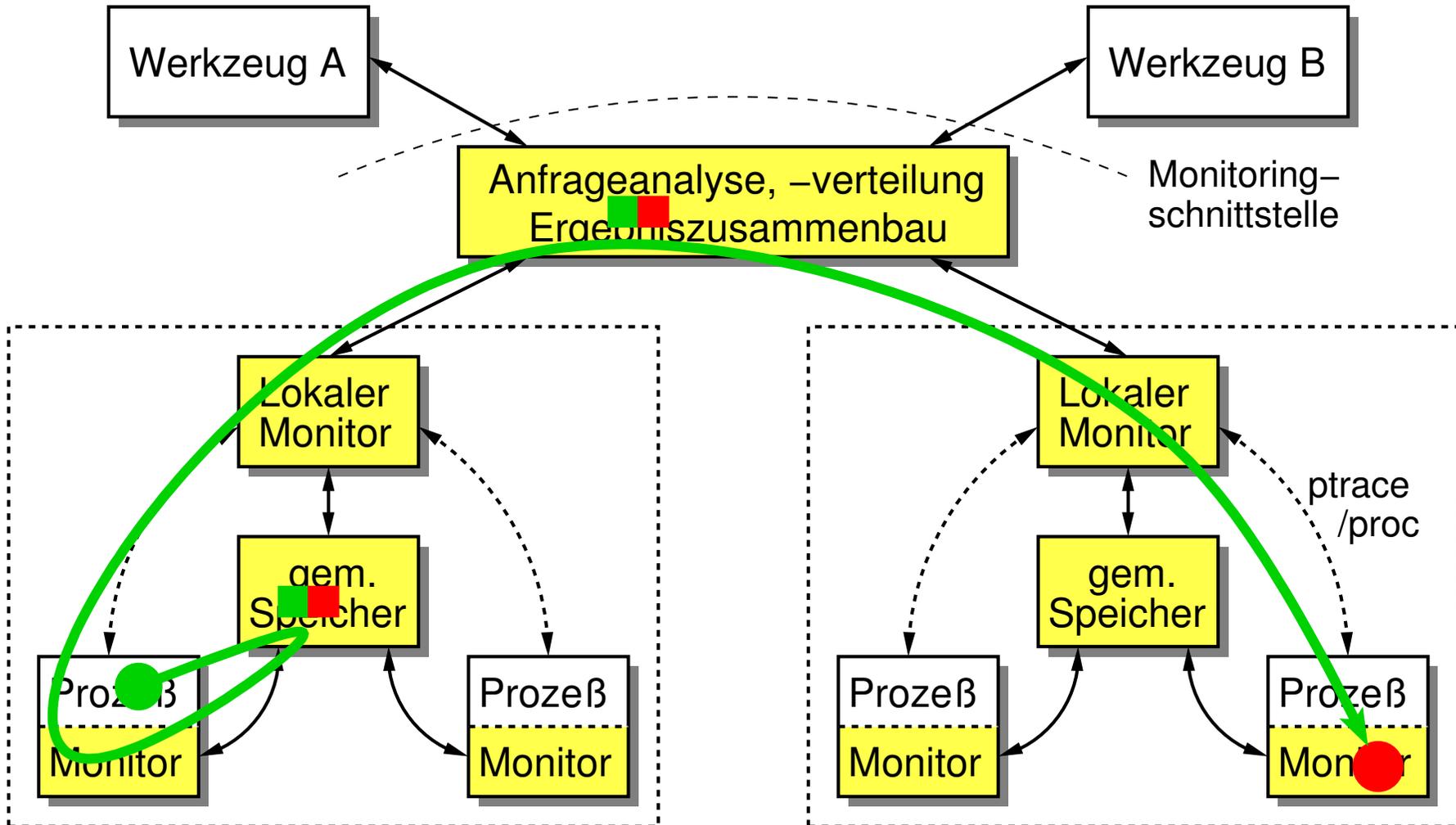


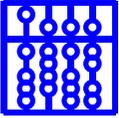
Mediatoren / Interaktoren

- ➔ Mediator ist zusätzliche Komponente
- ➔ Überwacht Ereignisse, führt zugehörige Aktionen aus



Anfragebearbeitung im OCM





Spezielle Eigenschaften des OCM:

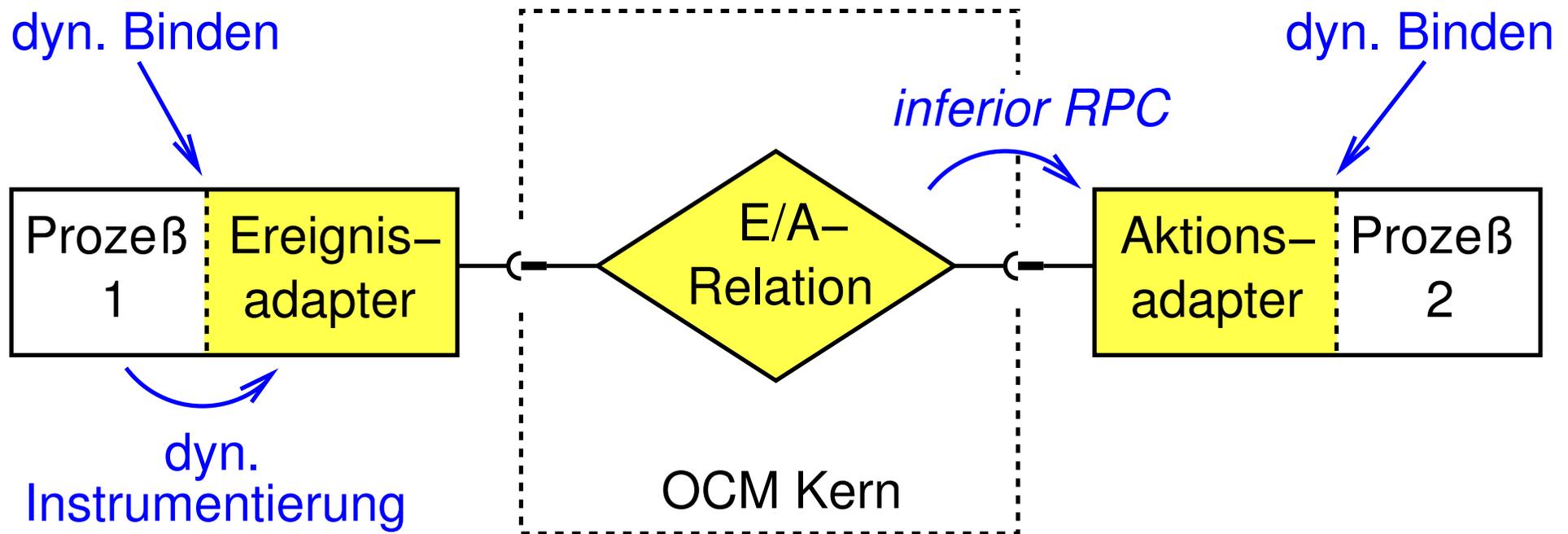
- ➔ Überwachung von Ereignissen in Prozessen
 - ➔ Aufruf von Aktionen in Prozessen
 - ➔ Menge der Ereignisse und Aktionen zur Laufzeit erweiterbar
- } ohne spezielle Vorbereitung dieser Prozesse

Löst wichtiges Problem:

- ➔ Komponenten müssen a-priori Schnittstellen bereitstellen
 - ➔ Vorhandensein?
 - ➔ Vollständigkeit / Angemessenheit?
 - ➔ Abhängigkeit von Middleware

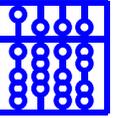
Zur Realisierung:

- ➔ Kombination von dynamischer Instrumentierung, dynamischem Binden und *inferior RPC*



Verallgemeinerung: Monitorsystem = dynamische Middleware

- ➔ Mediatorkonzept
 - ➔ Verhaltensregeln dynamisch programmierbar
- ➔ Dynamisches Anbinden von Schnittstellen(adapttern)
 - ➔ alle für Binder sichtbaren Symbole nutzbar
 - ➔ globale Variable / Objekte
 - ➔ globale Funktionen
 - ➔ öffentliche und geschützte Methoden
 - ➔ Ereignisse: Beginn / Ende einer Funktion bzw. Methode
 - ➔ Aktionen: nutzen Komponente wie Programmierbibliothek



Konsequenzen:

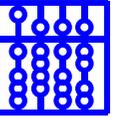
- ➔ Flexiblere Nutzung von Komponenten
 - ➔ Mächtigere (interne) Schnittstelle nutzbar
 - ➔ Einfache und effiziente Anpassung der Schnittstelle

- ➔ Vereinfachte Entwicklung von Komponenten
 - ➔ Kein expliziter Export von Schnittstellen
 - ➔ kein unnötiger 'Ballast' im Code
 - ➔ keine Abhängigkeit von einer Middleware

- ☞ Schnittstelle an sich / Dokumentation muß vorhanden sein!

Prototyp der Werkzeugumgebung

- ➔ Sichern des Programmzustands, automatisches Rücksetzen
- ➔ Konsistenz VISTOP ↔ DETOP
- ➔ Zur Verhaltensintegration: 3 Ereignisse + 10 Aktionen
 - ➔ nur eine Aktion mit komplexerem Adapter
 - ➔ ohne Modifikation des Werkzeug-Quellcodes
- ➔ Aber: Datenintegration erforderte Modifikationen (Konflikte durch Verletzung impliziter Annahmen)
- ☞ Werkzeuge müssen Vorhandensein anderer Werkzeuge prinzipiell berücksichtigen, nicht aber konkrete Kooperation



Laufzeit-Werkzeuge

- Beispiele, Forschungsthemen

Monitorsysteme

- Universelle Schnittstelle
- Effiziente Implementierung

Interoperable Laufzeit-Werkzeuge

- Motivation, Beispiel
- Datenintegration
- Verhaltensintegration

Künftige Arbeiten

- ➔ Automatische Leistungsanalyse interaktiver Anwendungen in Grid-Umgebungen (Esprit-Projekt CrossGrid)
- ➔ Interaktion zwischen Werkzeugen und Compilern
 - ➔ Debugging / Leistungsanalyse von HPF+ Programmen
 - ➔ Adaptive Laufzeitsysteme
- ➔ Konsolidierung des OCM, weitere Zielplattformen
- ➔ Ausbau der Konzepte zur Verhaltensintegration
 - ➔ Übertragung in andere Anwendungsbereiche
 - ➔ Integration mit Koordinationssprache