

Analiza wydajnościowa klastrowego systemu internetowego

Tomasz RAK

Katedra Informatyki i Automatyki
Politechnika Rzeszowska

6 maja 2015

Klastry...



Agenda

- Modelowanie
- Rozproszony system webowy
- Symulacje z użyciem QPME

Wcześniejsze prace

Modelowanie z użyciem QPN ^a

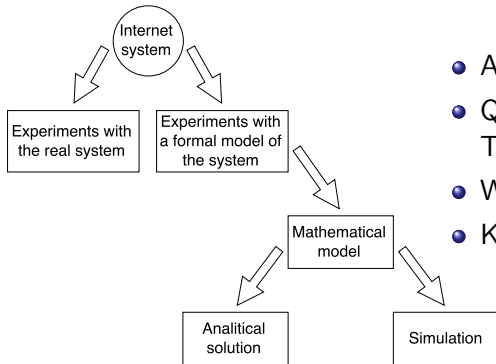
^aRak T.: **Performance Analysis of Distributed Internet System Models Using QPN Simulation**, Computer Science and Information Systems (FedCSIS) (2014)
doi:10.15439/2014F366

Rak T.: **Performance Analysis of Cluster-Based Web System Using the QPN Models**, Information Sciences and Systems, Springer, pp. 239-247 (2014)
doi:10.1007/978-3-319-09465-6_25

Agenda

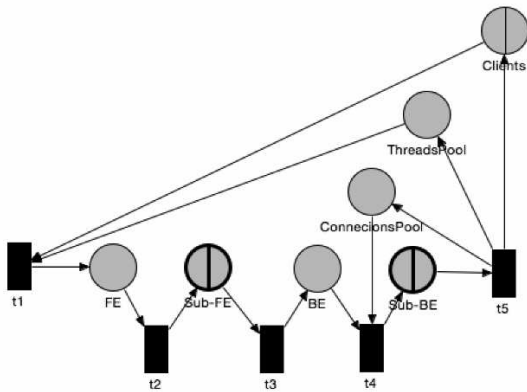
- Modelowanie - teoria
- Rozproszony system webowy \leftrightarrow 12
- Symulacje z użyciem QPME \leftrightarrow 14

Rozwiązania wydajnościowe

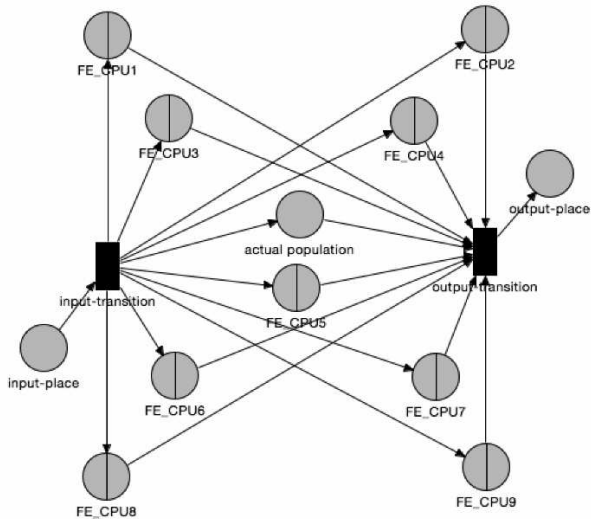


- Analiza wydajnościowa
- $QPN = PN + QN$ (wcześniej TCPN, QN)
- Warstwy (FE, BE)
- Klaster

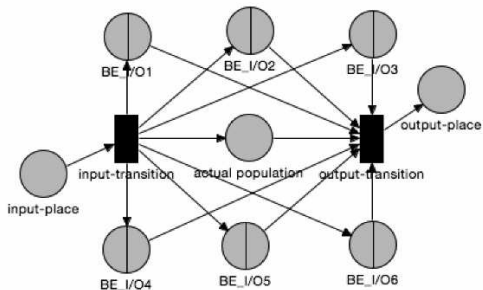
Model DWS



Model DWS (FE)



Model DWS (BE)



$$CPN = (P, T, C, I, M) \quad (1)$$

where:

- $P = \{p_1, p_2, \dots, p_i\}$ is a finite and non-empty set of places,
- $T = \{t_1, t_2, \dots, t_j\}$ is a finite and non-empty set of transitions,
- $P \cap T = \emptyset$,
- C is a colour function defined from $P \cup T$ into finite and non-empty sets (specify the types of tokens that can reside in the place and allow transitions to fire in different modes),
- $I(p, t)$ are the backward and forward incidence functions defined on $P \times T$, such that $I(p, t) \in [C(t) \rightarrow C(p)]$, $\forall (p, t) \in P \times T$ (specify the interconnections between places and transitions),
- $M(p)$ is a initial marking defined on P such that $M(p) \in C(p)$, $\forall p \in P$ (specify how many tokens are contained in each place).

$$QPN = (CPN, Q, W) \quad (2)$$

where:

- $Q = (Q_1, Q_2, (q_1, \dots, q_{|P|}))$, where:
 - $Q_1 \subseteq P$ is a set of timed queueing places,
 - $Q_2 \subseteq P$ is a set of immediate queueing places,
 - $Q_1 \cap Q_2 = \emptyset$,
 - $(q_1, \dots, q_{|P|})$ is an array with description of places (if p_i is a queueing place q_i denotes the description of a queue with all colors of $C(p_i)$ into consideration or if p_i is the ordinary place (p_i) equals *null*).
- $W = (W_1, W_2, (w_1, \dots, w_{|T|}))$, where:
 - $W_1 \subseteq T$ is a set of timed transitions,
 - $W_2 \subseteq T$ is a set of immediate transitions,
 - $W_1 \cap W_2 = \emptyset$, $W_1 \cup W_2 = T$,
 - $(w_1, \dots, w_{|T|})$ is an array (entry $w_j \in [C(t_j) \mapsto \mathbb{R}^+]$ such that $\forall c \in C(t_j) : w_j(c) \in \mathbb{R}^+$) of:
 - rate of a negative exponential distribution specifying the firing delay due to colour, if $t_j \in W_1$,
 - firing weight specifying the relative firing frequency due to colour, if $t_j \in W_2$.

$$QPN = (P, T, C, I, M, Q, W) \quad (3)$$

where:

- $P = \{FE, BE, ThreadsPool, ConnectionsPool\}$,
- $T = \{t_1, t_2, t_3, t_4, t_5\}$,
- $C(t_j)$, where j is number of transition,
- $I(p, t)$,
- $M(p)$,
- $Q = (Q_1, Q_2, (-/M/\infty/IS/\infty_{Clients}, null, -/M/1/PS/\infty_{Sub-FE}, null, -/M/1/FIFO/\infty_{Sub-BE}, null, null))$,

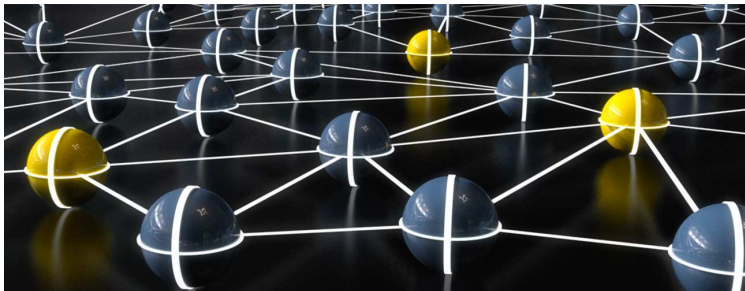
where:

- $Q_1 = \{Clients, FE_CPU_m, BE_I/O_n\}$,
- $Q_2 = \emptyset$,
- $W = (W_1, W_2)$, where:
 - $W_1 = \emptyset$,
 - $W_2 = T$,
 - $\forall c \in C(t_j) : w_j(c) := 1$ (all transition firings are equally likely).

Agenda

- Modelowanie ↔ 5
- Rozproszony system webowy - DWS
- Symulacje z użyciem QPME ↔ 14

Klastry



Agenda

- Modelowanie \rightarrow 5
- Rozproszony system webowy \rightarrow 12
- Symulacje z użyciem QPME - analizy

Rodzaje symulacji

Wydajność systemu w zależności od:

- obciążenia
- rozmiaru puli wątków FE i połączeń BE
- zmiennej liczby elementów klastra

Parametry, od których zależy czas odpowiedzi

- *Service Demand, Residence Time*
- *Workload Intensity*

Czas odpowiedzi (response time) jest równy sumie czasów obsługi w poszczególnych zasobach (residence time), gdzie: i - liczba miejsc:

$$R = \sum_i^{k=1} R'_k \quad (4)$$

Czas obsługi w zasobie (residence time) jest sumą czasu spędzonego w kolejce (queueing time) i średniego czasu obsługi dla zasobu (service demand):

$$R'_k = Q_k + D_k \quad (5)$$

gdzie, czas spędzony w kolejce (czas oczekiwania) na zasób to $Q_k = \sum_i^{k=1} q_k$ i średni czas obsługi w określonym zasobie to $D_k = \sum_i^{k=1} d_k$.

Średni czas obsługi w określonym zasobie, z wyłączeniem czasu oczekiwania na zasób. Nie zależy od obciążenia!

Całkowity czas odpowiedzi

Całkowity czas odpowiedzi:

$$\begin{aligned} R_{TOTAL} = & R_{CLIENT_{(DEPOSITORY)}} + R_{FE_{(QUEUE)}} + \\ & \sum_{i=1}^n R_{FE_CPUi_{(QUEUE)}} + \sum_{i=1}^n R_{FE_CPUi_{(DEPOSITORY)}} + \\ & R_{BE_{(QUEUE)}} + R_{BE_CPU_{(QUEUE)}} + R_{BE_IO_{(QUEUE)}} \end{aligned} \quad (6)$$

- $R_{FE_{(QUEUE)}}$ i $R_{BE_{(QUEUE)}}$ - są to miejsca użyte do zatrzymania przychodzących zapytań, gdy oczekują one na wątki (serwer aplikacji) i procesy (baza danych)

Parametry (1)

Średni czas obsługi dla węzłów w poszczególnych warstwach:

- $d_{FE_CPU} = 0,714$ [ms]
- $d_{BE_I/O} = 0,133$ [ms]

Znakowanie początkowe dla miejsc:

- liczba klientów (liczba znaczników w miejscu *Clients*)
- liczba wątków serwera FE (liczba znaczników w miejscu *ThreadsPool*)
- liczba połączeń do serwera BE (liczba znaczników w miejscu *ConnectionsPool*)

Typy znaczników w miejscu (kolor):

- zapytania (jedna klasa)
- wątki
- połączenia

Parametry (2)

	Elementy/Parametry	Nazwa/Wartość
QPME	Miejsce kolejkowe FE	FE_CPUm^a
	Miejsce kolejkowe BE	BE_I/On^b
Oprogramowanie ^c	Miejsce <i>ThreadsPool</i>	30^d
	Miejsce <i>ConnectionsPool</i>	40^e
Obciążenie	λ	0,015; 0,03; 0,045; $0,06^f$
	Miejsce kolejkowe <i>Clients</i> ^g	100, 200, 300, 400, 500
	Czas symulacji [s]	300

^a m - Liczba węzłów FE

^b n - Liczba węzłów BE

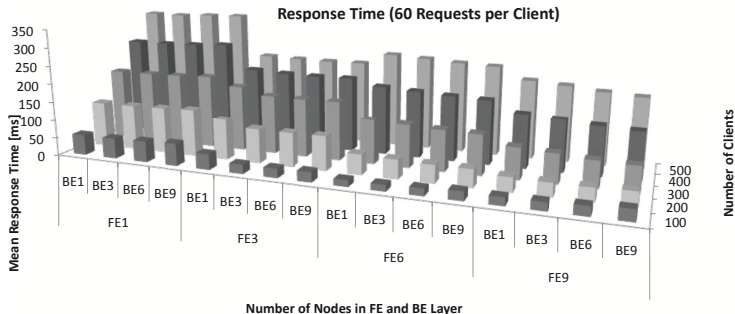
^c Znakowanie początkowe

^d 30 wątków dla jednego węzła FE, 90 wątków dla trzech węzłów FE, 180 wątków dla sześciu węzłów FE, 270 wątków dla dziewięciu węzłów FE

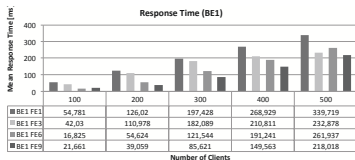
^e 40 połączeń dla jednego węzła BE, 120 połączeń dla trzech węzłów BE, 240 połączeń dla sześciu węzłów BE, 360 połączeń dla dziewięciu węzłów BE

^f Czas namysłu klienta 66,67; 33,33; 22,22; 16,67 [ms] przy obciążeniu: 15, 30, 45, 60 [zapytań na sekundę]

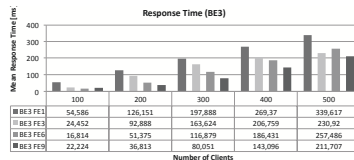
FE (1, 3, 6 i 9) i BE (1, 3, 6 i 9) oraz różna liczba klientów (6000-30000)



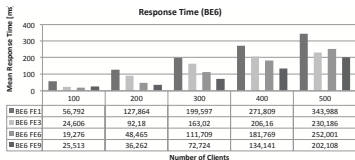
BE (1, 3, 6 i 9) i FE (1, 3, 6 i 9) oraz różna liczba klientów (6000-30000)



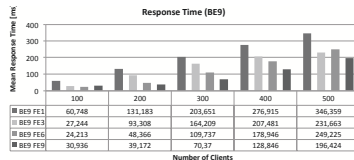
a)



b)



c)



d)

Wnioski z przeprowadzonych prac

- Przygotowano modele symulacyjne
- Analiza wyników

Przyszłe prace

Charakterystyka obciążenia (model klienta):

- Klasy zapytań
- Użycie elementów sprzętowo-programowych przez poszczególne klasy

Różne typy/klasz klientów to różnz zachowania

Różnz sposoby przybywania zapytań, popytu na usługi czy trasowania sieciowego.

Response Time Analysis of Distributed Web Systems Using QPNs

Mathematical Problems in Engineering (Mathematical Problems in Petri Nets Theory and Applications)

Hindawi Publishing Corporation

Dziękuję za uwagę!

Modelowanie 5

Rozproszony system webowy 12

Symulacje z użyciem QPME 14

Simulations can be used as experiments.