

# VU Research Portal

## Performance models for analysis and control of it systems

Hristov, A.

2018

### **document version**

Publisher's PDF, also known as Version of record

[Link to publication in VU Research Portal](#)

### **citation for published version (APA)**

Hristov, A. (2018). *Performance models for analysis and control of it systems*. [PhD-Thesis - Research and graduation internal, Vrije Universiteit Amsterdam].

### **General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

### **Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

### **E-mail address:**

[vuresearchportal.ub@vu.nl](mailto:vuresearchportal.ub@vu.nl)

## Samenvatting

De toenemende complexiteit van IT infrastructures vormt een grote uitdaging voor de beheersbaarheid van computersystemen. Hierbij is het van groot belang de beschikbare bronnen efficiënt in te zetten onder voorwaarde dat de gewenste Quality of Service (QoS) wordt behaald. Om deze complexiteit en schaal het hoofd te kunnen bieden moet er verder worden gedacht dan voor de hand liggende ad-hoc oplossingen. Dit heeft ons geïnspireerd om prestatie modellen en algoritmen te ontwikkelen die het gedrag van ICT-serviceketens beschrijven. De modellen en algoritmen kunnen bovendien gebruikt worden voor het bepalen van besturingsregels voor deze ICT-systemen. Een van de hoofddoelen van ons onderzoek is het overbruggen van de kloof tussen theorie en praktijk. Om dit te bewerkstelligen plaatsen we elk van de bestudeerde modellen en algoritmen in de context van een beoogde ICT-toepassing. We introduceren een nieuwe oplossingsmethode waarin technieken uit de wachtrijanalyse en machine learning vakgebieden worden gecombineerd. In deze methode passen we machine learning toe op concepten uit de wachtrijanalyse. De resulterende modellen zijn in staat om variëteit van relevante externe invloeden mee te nemen. Dit maakt de ontwikkeling mogelijk van accurate en schaalbare modellen en algoritmen met een bredere toepasbaarheid.

Het model in hoofdstuk 1 is geïnspireerd door de interactie tussen database toepassingen en de onderliggende opslag. Deze interactie vindt plaats doormiddel van een cache die schrijfverzoeken tijdelijk opslaat in werkgeheugen. Ons model kan worden gezien als een tandemmodel waarbij de uitvoer van de eerste server genereert de invoer van de tweede server. De eerste server representeert het verzenden van schrijfverzoeken vanuit applicatieniveau naar de cache. De tweede server modelleert het schrijven van de cache naar de onderliggende opslag. Elk van de twee servers kan op elk gewenst moment worden in- of uitgeschakeld. Daarnaast veronderstellen we wachtkosten bij elk van de twee wachtrijen en opstartkosten voor de tweede server, d.w.z. de cache. We construeren een Markov Decision Process (MDP) om de optimale strategie te bepalen die de gemiddelde langetermijnkosten van het systeem minimaliseert. De bestaande numerieke technieken om de bedieningsstrategie te optimalis-

eren zijn inefficiënt en zelfs onhaalbaar voor bepaalde systemen in de realiteit. Daarom presenteren we een schaalbare benadering om de dempelwaardestrategie te bepalen voor het in en uitschakelen van de servers. De onderzochte techniek is intuïtief en resulterende strategie zorgt dat de gemiddelde kosten over de lange termijn binnen enkele procenten van optimale strategie liggen.

Hoofdstuk 3 beschouwt de prestatie modellen en bijhorende optimalisering van database caching-mechanismen vanuit een ander perspectief. In tegenstelling tot het vorige hoofdstuk, richt dit hoofdstuk zich op het zogenaamde ‘write-behind’ cachemechanisme. We modelleren het systeem doormiddel van een Markovmodel waarin een enkele server de betreffende taken in twee stappen verwerkt. De eerste stap correspondeert met plaatsen en voorbereiden van de binnenkomende taken in de cache. Na de voorverwerking worden de taken verzameld in een batch ter grootte  $K$  die in één keer wordt verwerkt als een grote schrijfofdracht. Het doel is om de grootte van de batch te bepalen die de verwachte wachttijd voor aankomende opdrachten minimaliseert. Om de optimale waarde van  $K$  te schatten, leiden we eerst een analytische oplossing af voor de stationaire verdeling van taken in het systeem voor twee speciale gevallen:  $K = 1$  en  $K = 2$ . Vervolgens leiden we een vloeistofmodel af om de verwachte wachttijd van grote batches te benaderen. Ten slotte combineren we deze technieken tot een uitdrukking waarmee de optimale batchgrootte op een efficiënte en schaalbare manier kan worden bepaald. Uitgebreide numerieke experimenten tonen aan dat de benadering heel goed werkt voor een veelvoud van parametercombinaties.

In hoofdstuk 4 introduceren we een nieuwe techniek om expliciete benaderingen te bepalen voor drempelwaardestrategieën in MDP’s. Onze methode maakt gebruik van het Symbolic Regression (SR) -algoritme. We zijn in staat zowel de nauwkeurigheid als de rekentijd van dit algoritme aanzienlijk verbeteren door gebruik te maken van domeinspecifieke kennis uit het MDP-raamwerk. We passen onze aanpak toe op twee voorbeelden. Voor deze voorbeelden hebben we analytische uitdrukkingen afgeleid die de optimale strategie met grote nauwkeurigheid benaderen. De gesloten vorm van deze uitdrukkingen stelt ons bovendien in staat een gevoeligheidsanalyse uit te voeren op de systeemp parameters. De gevonden uitdrukkingen maken het mogelijk om drempelwaardestrategieën direct aan te passen aan veranderingen in de systeemp parameters. Dit maakt onze oplossing bijzonder aantrekkelijk voor een toepassing in real-time systemen. We zijn er van overtuigd dat onze aanpak zeer generiek is en van toepassing is op een breed spectrum van MDP formuleringen.

In hoofdstuk 5 breiden we ons onderzoek uit naar de combinatie van wachtrijtheorie met technieken uit de machine learning. We introduceren een manier om de technieken en modelcomponenten uit wachtrijtheorie te vertalen naar SR. De gepresenteerde methode kan gesloten vorm benaderingen afleiden voor de belangrijkste prestatie-maten van wachtrijmodellen. Ter demonstratie passen we onze aanpak toe op het zogenaamde ‘two-stream blending’ systeem. Ook voor dit model hebben we uitdrukkingen in gesloten vorm afgeleid voor de gemiddelde wachttijd en de maximale doorvoersnelheid. De gegenereerde analytische uitdrukkingen zijn nauwkeurig en tegelijkertijd algebraïsch eenvoudig.

Tenslotte beschouwen we een model dat op een hoger abstractieniveau ligt dan de modellen in voorgaande hoofdstukken. In hoofdstuk 6 analyseren we een Layered Queuing Network (LQN) waarin servers een gelaagde keten vormen en daarom elkaars doorvoersnelheid beïnvloeden. We presenteren een eenvoudig benaderingsalgoritme voor het bepalen van de bezettingsgraad en doorvoersnelheid van de gelaagde servers. Wij geloven dat de hoge nauwkeurigheid van de benadering potentie biedt voor uitbreiding naar generiekere gelaagde wachtrijmodellen. Verder laten we zien dat zelfs wanneer prestatie-maten zoals de maximale doorvoersnelheid en de bezettingsgraad van de servers bekend zijn, het bepalen van de beperkende factor in het netwerk niet triviaal is. Daarom introduceren we een prestatie-index die een uitbreiding is van de ‘slowest server rule’. Uit de uitgevoerde numerieke tests blijkt dat de voorgestelde bottleneck-identificatie techniek nauwkeuriger is dan de bestaande algoritmen. Daarnaast is onze methode ook toepasbaar op klassieke wachtrijnetwerken zonder gelaagde structuur.