

# Causal study of Network Performance

Hadrien Hours<sup>1</sup>, Ernst Biersack<sup>1</sup> and Patrick Loiseau<sup>1</sup>

<sup>1</sup>*EURECOM, Campus SophiaTech, Les Templiers, 450 Route des Chappes, 06410 Biot Sophia Antipolis, France*

---

The use of Internet in the every day life has pushed its evolution in a very fast way. The heterogeneity of the equipments supporting its networks, as well as the different devices from which it can be accessed, have participated in increasing the complexity of understanding its global behavior and performance. In our study we propose a new method for studying the performance of TCP protocol based on causal graphs. Causal graphs offer models easy to interpret and use. They highlight the structural model of the system they represent and give us access to the causal dependences between the different parameters of the system. One of the major contribution of causal graphs is their ability to predict the effects of an intervention from observations made before this intervention.

**Keywords:** Bayesian graph, Causality, Network Performance, TCP Protocol

---

## 1 Introduction

The TCP protocol supports more than 80% of the traffic going through the Internet and several studies of its behavior exist to better understand the different parameters and their roles in the performance one can expect when using this protocol. Theoretical models [Padhye 98] as well as statistical ones [Mirza 07] have been proposed. These models do not take into account the impact that parameters can have between each others and the approach is highly depending on the TCP version that is considered. For the second class of models, as relying on statistical correlation, there is the additional limitation of creating a model that is invalid as soon as it is used for predicting intervention that could change the statistics under which the model was inferred. The causal approach we propose in this paper answers these different issues.

The different works made in the domain of causal models inference and their representations as Bayesian graph, often Directed Acyclic Graphs (DAGs), as support for predicting interventions [Pearl 09, Spirtes 01] give access to new perspectives in the study and understanding of complex systems.

While causality has been used to study network performance [Tariq 08], we place ourselves in a situation where resource, both in terms of data and computational power, is a limiting factor. We also introduce what we believe to be more formal and robust methods for both graph inference and intervention predictions.

In this paper we present the Causal Model Inference approach in Section 2 and show one example of its application to the study of FTP traffic in Section 3. We will conclude this paper in Section 4.

## 2 Background

Correlation is not causality, two parameters can be found correlated but this correlation does not give us any information on whether one is the cause of the other, or the opposite, or if there exists a third parameter *causing* these two parameters. This basic notion illustrates the big difference between causal model and statistical model.

### 2.1 Causal model inference

In our work we use the PC algorithm [Spirtes 01] to infer our causal graph. This algorithm is based on leading independence tests to, gradually, build a graph where these independences are verified. The key, then, is the choice of the independence test that will assess the properties of our system. While the classical approach is to test partial correlation between the residuals of linear regression (as in Z-Fisher test), we are

Parameter	Definition	Min	Max	Avg	Coeff. Var.
Dist (km)	Distance between Server and Client	14	620	250	0.95
T.O.D. (s)	Time of the day, when the connection was started	740	82000	46000	0.53
Nbbytes (MB)	Number of bytes sent by the server	6	60	31.5	0.51
N.L.A.C. (kbps)	Narrow Link available capacity	47.8	42.9e+3	5.9e+3	1.7
Nbhops (units)	Number of hops between Client and Server	9	27	11	0.24
RTT (ms)	Round Trip Time	60	710	270	0.76
BufferingDelay (ms)	Part of the RTT due to queuing delay	4.2	470	84	1.2
RetrScore (no unit)	Fraction of retransmitted packets	0.001	0.014	0.0018	0.92
p (no unit)	Fraction of retransmitted windows of packets	3e-5	0.011	0.0007	1.3
T.O.R. (no unit)	Fraction of retransmitted packets due to Time Outs	0	0.01	0.0006	1.7
RWin (kbytes)	Receiver Window	10.7	253.9	136.7	0.68
Tput (kbps)	Throughput	24	928	332	0.77

Table 1: Summary of FTP traffic dataset

observing parameters which are not normally distributed and with non linear dependences. Consequently, in our work, we use the Hilbert Schmidt Independence Criterion [Zhang 12], a criterion that does not rely on normality or linearity to test parameters independences.

## 2.2 Prediction

Assuming now that we have our causal graph representing our system, it is possible, using simple graphical criteria [Pearl 09, Spirtes 01], to predict the effects of an intervention on one of the parameters of our model.

In the graph each vertex corresponds to one parameter of our system and an edge between two vertices represents a direct causal effect from the parent to the child (in  $X \rightarrow Y$ ,  $X$  is the parent and  $Y$  the child). A (set of) node(s),  $Z$ , is said to *block* a path between  $X$  and  $Y$  if (i) every collider ( $\rightarrow W \leftarrow$ ), or its descendants, is not in  $Z$  and (ii) at least one non collider is in  $Z$ .

The notation  $do(X = x)$  represents the manual intervention of setting the parameter  $X$  to the value  $x$ .

**Definition 1 (Back-door criterion)** A set of variables,  $Z$ , is said to satisfy the Back-door criterion, relative to an ordered pair of variables  $(X_i, X_j)$ , in a DAG  $G$  if: (i) No node in  $Z$  is a descendant of  $X_i$ . (ii)  $Z$  blocks every path between  $X_i$  and  $X_j$  that contains an arrow into  $X_i$

**Theorem 1 (Back-door adjustment)** If a set of variables  $Z$  satisfies the Back-door criterion relative to  $(X, Y)$ , then the causal effect of  $X$  on  $Y$  is identifiable and given by the formula:

$$P(Y = y | do(X = x)) = \sum_Z P(Y = y | X = x, Z = z)P(Z = z) \quad (1)$$

Due to space constraints we will not present the theory of *d-separation*, at the source of this criterion, but we redirect the reader to [Pearl 09] for a formal explanation and justification of the Back-door criterion.

## 3 Study of FTP traffic

For illustrating our approach we decided to study FTP traffic. One reason for choosing this protocol is the absence of the application limiting the performance of TCP. The throughput is only limited by the Network or the client Receiver Window.

We first define the different parameters of our model and explain how do we build our dataset. Then, we present the causal model we obtain, with the PC algorithm, before predicting the effect of intervening on the *RTT*, on the *Throughput*.

### 3.1 Dataset

For our experiments we set up a FTP server where all the traffic is recorded. Using Intrabase [Siekkinen 05] and the Tstat tool<sup>†</sup> we obtain the different metrics presented in Table 1. This table also presents a summary of our dataset consisting in 1000 downloads from different clients, in Spain, Germany and France.

<sup>†</sup> <http://tstat.tlc.polito.it/>

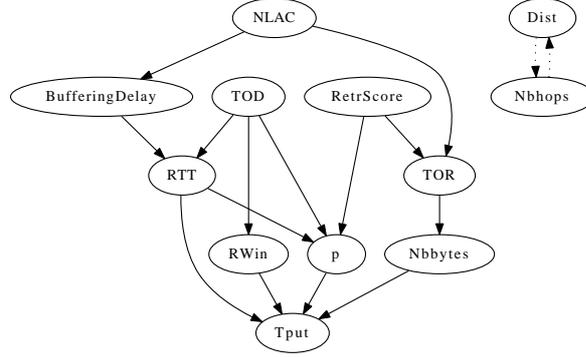


Figure 1: Causal graph model of FTP traffic performance

### 3.2 Causal model

The model we obtained is presented in Figure 1. Due to space constraints we cannot describe all the properties of this model and focus on prediction, presented in the following section.

### 3.3 Predictions

As we are working with variables where no distributions can be assumed, we use Copula Theory [Jaworski 10] to estimate the densities, and conditional densities, of the parameters present in the Back Door adjustment formula, Equation (1).

The presentation of Copulae, as well as the choice of the Copula to model the distributions of our parameters, are out of the scope of this paper.

For predicting the effect of intervening on  $RTT$ , we use the Back Door Criterion met by the set of variables  $Z_{RTT}$ , Equation (2), with  $Z_{RTT} = \{T.O.D., N.L.A.C.\}$ :

$$f_{TPUT|do(RTT)}(Tput = \Delta | do(RTT = \theta)) = \int_{Z_{RTT}} f_{Tput|RTT,Z_{RTT}}(Tput = \Delta | RTT = \theta, Z_{RTT}) f_{Z_{RTT}} dZ_{RTT} \quad (2)$$

Figure 2 presents the estimated *Throughput* post-intervention (solid line), obtained with Equation (2). For comparison, we present the *Throughput* obtained from observations, in the initial dataset, for which the  $RTT$  value is the one corresponding to the intervention we want to predict (dotted line). We also added the information of the number of samples (bar plot) from which the distribution of the *Throughput* post-intervention was estimated. For values of  $RTT$  in  $[250, 350]$  we can see that there are not enough samples for estimating the post-intervention distribution which leads to inconsistent estimates.

Removing the inconsistent estimates of the post-intervention distribution, we can see that the *Throughput* post intervention is smaller than the one we observe in the initial dataset for a given value of  $RTT$ . This observation can be explained by the causal graph presented Figure 1 as, by conditioning on a given  $RTT$  value, we also take into account the Back door effects of variables such as  $RWin$ ,  $N.L.A.C.$ ,  $T.O.D.$  or  $T.O.R.$  which are spurious associations blocked by  $Z_{RTT}$  in Equation (2). This result shows that adopting a naive approach of estimating the *Throughput* directly from the pre intervention samples will overestimate the effect that an intervention on the  $RTT$  will have on the *Throughput*.

## 4 Conclusion

This work is a first attempt to apply causal theory to the study of network performance. We present the inference of causal models, represented by Directed Acyclic Graphs, and their use to predict the effects of interventions on complex systems from passive observations. We show the different challenges that arise when applying theoretical theory to a real case study, the study of FTP traffic, and propose several solutions. As every model, causal models present their limits. We present the ones we judge as the most important for our work and propose solutions to overcome them. We introduce several methods that show very promising

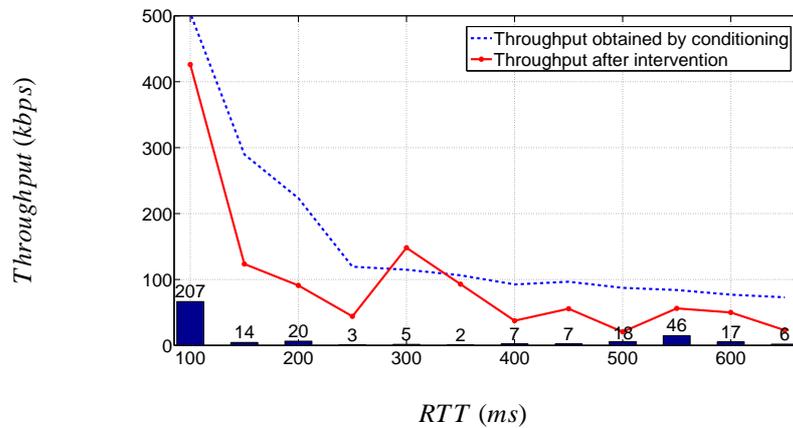


Figure 2: Comparison between causal approach and naive approach

results. The use of the Hilbert Schmidt Criterion, for testing independences in the PC algorithm, and copulae, for estimating the multidimensional probability density functions, are the two main ones.

As we can see from Equation (1), the variety of situations observed defines the range of predictions that we can make. It will be necessary to increase the number of samples we have to reach a greater accuracy in our predictions. We are now working on network simulated experiments where we will have access to more parameters and the possibility to verify the accuracy of our predictions.

One important limitation of our approach is the definition of static parameters, by averaging metrics such as  $RTT$  or  $RWin$ , to model a dynamic protocol. Using the web10G tool, we plan to have access to low level TCP parameters and to sample the parameters at a finer grained timescale.

## References

- [Jaworski 10] P. Jaworski, F. Durante, W. Härdle and T. Rychlik, *Copula Theory and Its Applications*, Lecture Notes in Statistics, 2010.
- [Mirza 07] M. Mirza, J. Sommers, P. Barford and X. Zhu, “A machine learning approach to TCP throughput prediction”, SIGMETRICS ’07, pp. 97–108, New York, NY, USA, 2007, ACM.
- [Padhye 98] J. Padhye, V. Firoiu, D. Towsley and J. Kurose, “Modeling TCP throughput: a simple model and its empirical validation”, *SIGCOMM Comput. Commun. Rev.*, 28(4):303–314, October 1998.
- [Pearl 09] J. Pearl, *Causality: Models, Reasoning and Inference*, Cambridge University Press, New York, NY, USA, 2009.
- [Siekkinen 05] M. Siekkinen, E. W. Biersack, G. Urvoy-Keller, V. Goebel and T. Plagemann, “InTraBase: integrated traffic analysis based on a database management system.”, E. Al-Shaer, A. Pras and P. Owezarski, Eds., *E2EMON*, pp. 32–46, IEEE, 2005.
- [Spirtes 01] P. Spirtes, C. Glymour and R. Scheines, *Causation, Prediction, and Search*, The MIT Press, Cambridge, MA, USA, Second edition, January 2001.
- [Tariq 08] M. Tariq, A. Zeitoun, V. Valancius, N. Feamster and M. Ammar, “Answering what-if deployment and configuration questions with wise”, *SIGCOMM Comput. Commun. Rev.*, 38(4):99–110, August 2008.
- [Zhang 12] K. Zhang, J. Peters, D. Janzing and B. Schölkopf, “Kernel-based Conditional Independence Test and Application in Causal Discovery”, *CoRR*, abs/1202.3775, 2012.