

How to find a Needle in a Haystack ?

On the Detection of Anomalies in Large Traces

Jean-Marc.Vincent@imag.fr

Laboratoire d'Informatique de Grenoble, INRIA Team Polaris

ANR Marmote
June 2, 2016

Joint work with : Robin Lamarche-Perrin, Lucas Mello Schnorr, Damien Dosimont,
Guillaume Huard, and Yves Demazeau.
Work partially supported by ANR Geomedia, Songs and Marmote



HOW TO FIND A NEEDLE IN A HAYSTACK ?

- 1 THE PROBLEM : Extracting macroscopic information from microscopic measures**
- 2 ANOMALY : Heterogeneous (unexpected) Macroscopic Behavior**
 - SPACE : Localization Heterogeneity
 - TIME : Temporal Heterogeneity
 - SPACE/TIME : Behavior Heterogeneity
- 3 METHODOLOGY : Information based Aggregation**
 - MACROSCOPIC STATE : Information based approach
 - PARTITION : Algorithms and Complexity
- 4 SYNTHESIS AND OPEN QUESTIONS**

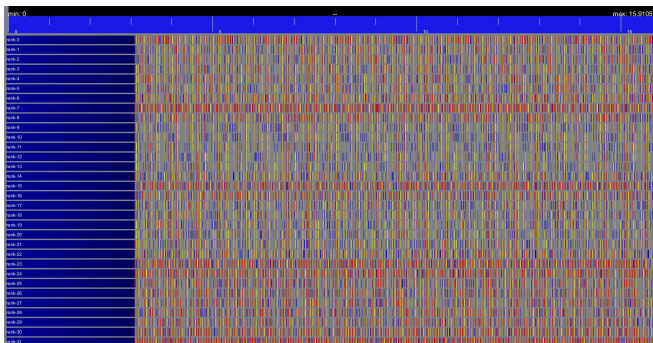
ARE WE ABLE TO UNDERSTAND THE BEHAVIOR OF LARGE APPLICATIONS ?



Stampede (TACC) ~ 500 000 cores

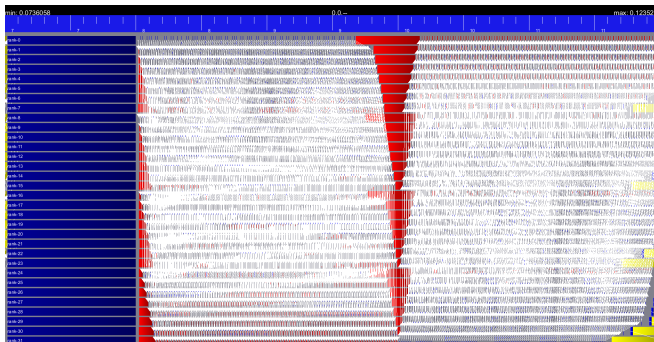
ANALYSIS OF LARGE TRACES : PERFORMANCE DEBUG

Adaptive analysis of large traces : example MPI program (LU factorization)
(Zoom in time)



ANALYSIS OF LARGE TRACES : PERFORMANCE DEBUG

Adaptive analysis of large traces : example MPI program (LU factorization)
(Zoom in time)



EXTRACTING MACROSCOPIC INFORMATION FROM MICROSCOPIC MEASURES

A haystack



Needles



Elementary model of processes behavior

Fine grain event traces or regular sampling

Id	date	location	state	event information
...
k	t_k	Core i	State 1	Event foo
$k+1$	t_{k+1}	Core i	State 2	Event bar
$k+2$	t_{k+2}	Core i	State 3	Event barfoo
...

Objective

- 1 Provide a measure of the quality of macroscopic visualizations.
- 2 Provide an interactive synthetic representation of large-scale data with partial multi-level aggregations.
- 3 Focus on heterogeneity with quantification

EXTRACTING MACROSCOPIC INFORMATION FROM MICROSCOPIC MEASURES

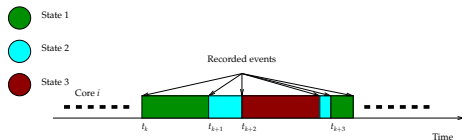
A haystack



Needles



Elementary model of processes behavior



Fine grain event traces or regular sampling

Id	date	location	state	event information
...
k	t_k	Core i	State 1	Event foo
$k + 1$	t_{k+1}	Core i	State 2	Event bar
$k + 2$	t_{k+2}	Core i	State 3	Event barfoo
...

Objective

- 1 Provide a measure of the quality of macroscopic visualizations.
- 2 Provide an interactive synthetic representation of large-scale data with partial multi-level aggregations.



EXTRACTING MACROSCOPIC INFORMATION FROM MICROSCOPIC MEASURES

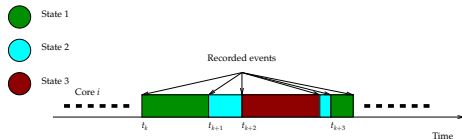
A haystack



Needles



Elementary model of processes behavior



Fine grain event traces or regular sampling

Id	date	location	state	event information
...
k	t_k	Core i	State 1	Event foo
$k + 1$	t_{k+1}	Core i	State 2	Event bar
$k + 2$	t_{k+2}	Core i	State 3	Event barfoo
...

Objective

- 1 Provide a measure of the quality of macroscopic visualizations.
- 2 Provide an interactive synthetic representation of large-scale data with partial multi-level aggregations.



HOW TO FIND A NEEDLE IN A HAYSTACK ?

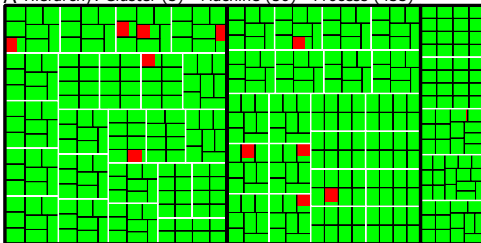
- 1 THE PROBLEM : Extracting macroscopic information from microscopic measures
- 2 **ANOMALY : Heterogeneous (unexpected) Macroscopic Behavior**
 - SPACE : Localization Heterogeneity
 - TIME : Temporal Heterogeneity
 - SPACE/TIME : Behavior Heterogeneity
- 3 METHODOLOGY : Information based Aggregation
 - MACROSCOPIC STATE : Information based approach
 - PARTITION : Algorithms and Complexity
- 4 SYNTHESIS AND OPEN QUESTIONS

SPACE ANOMALIES

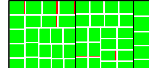
EXPERIMENTING ON WORK-STEALING TRACES

Case study 1 : classical behavior

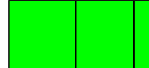
A Hierarchy: Cluster (3) - Machine (50) - Process (433)



A.1 Machine level



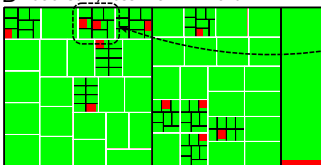
A.2 Cluster level



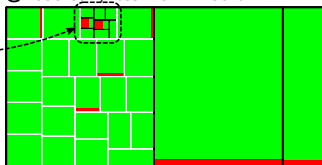
A.3 Full aggregation



B Ratio Gain/Loss with $P = 10\%$



C Ratio Gain/Loss with $P = 30\%$

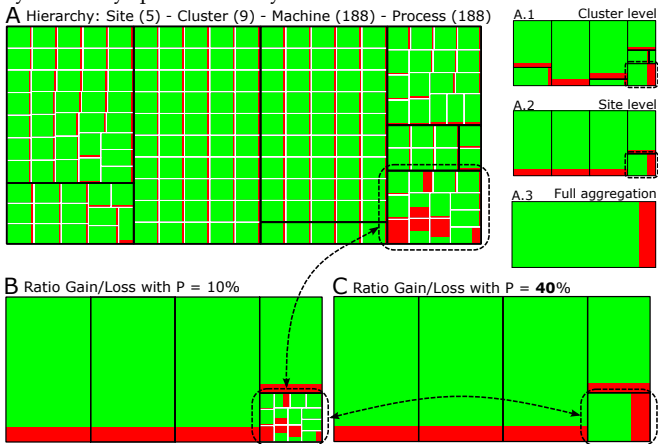


Multi-resolution representation: hierarchical partial zooming

SPACE ANOMALIES

EXPERIMENTING ON WORK-STEALING TRACES

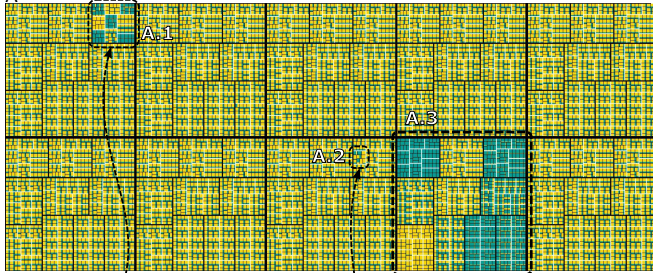
Case study 2 : widely spread anomaly



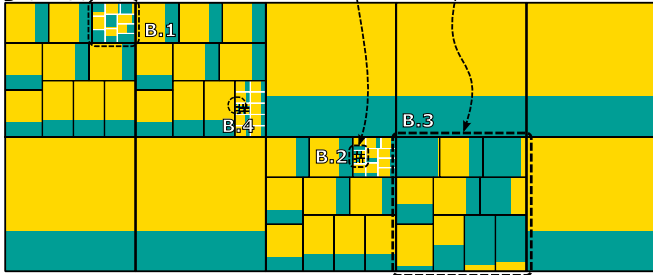
Multi-resolution representation : focus on heterogeneity

VISUALIZATION OF A MILLION OF PROCESSES

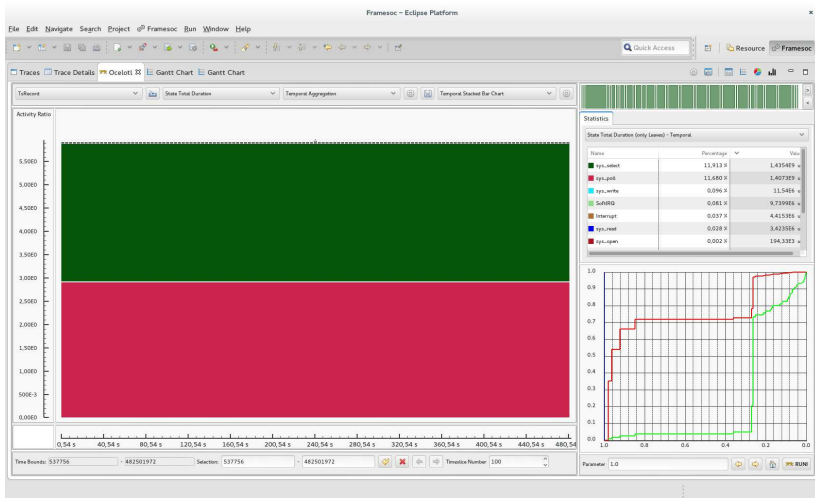
A Hierarchy: Site (10) - Super-Cluster (100) - Cluster (1000) - **Machine (10000)** - Process (1000000)



B with $P=10\%$



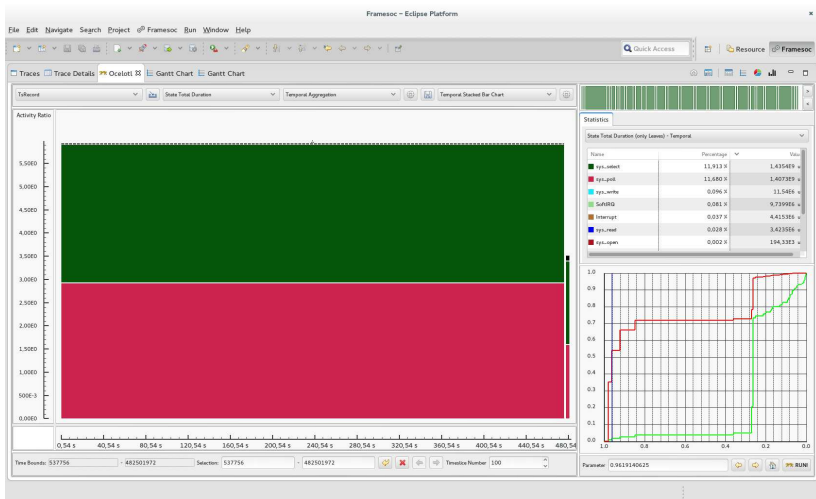
EXPERIMENTS ON EMBEDDED SYSTEMS



$p=1$

Multi-level aggregation: Ocelotl Application/Demo

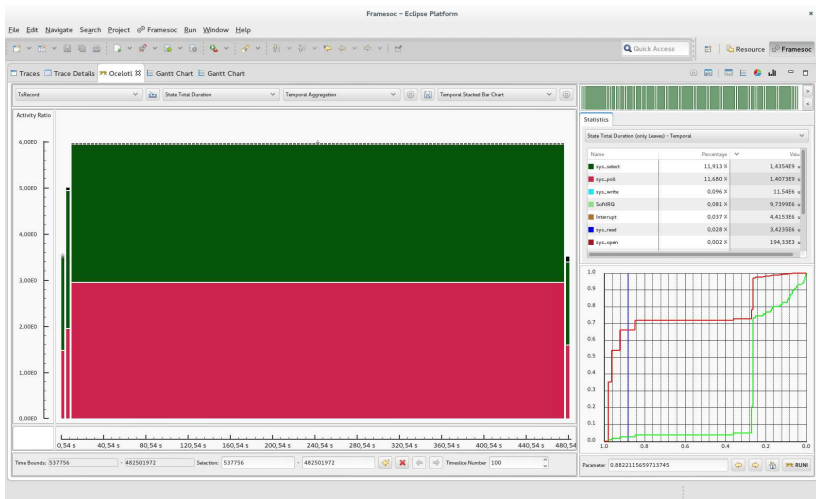
EXPERIMENTS ON EMBEDDED SYSTEMS



$p=0.96$

Multi-level aggregation: Ocelot Application/Demo

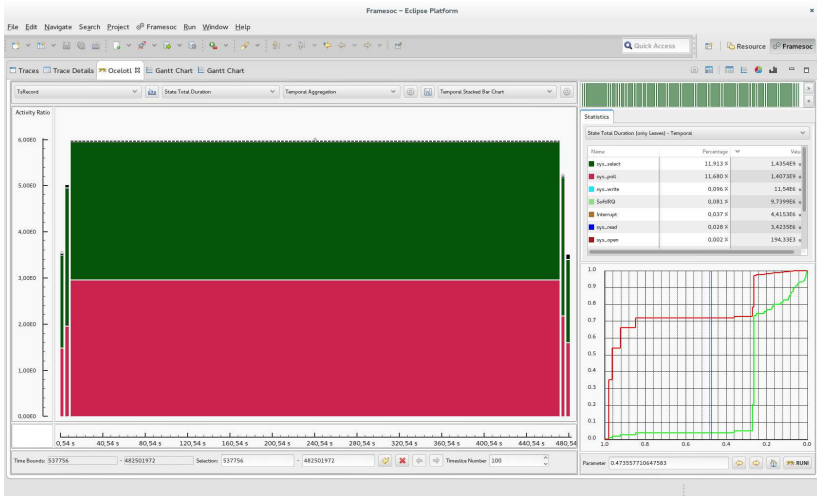
EXPERIMENTS ON EMBEDDED SYSTEMS



$p=0.88$

Multi-level aggregation: Ocelotl Application/Demo

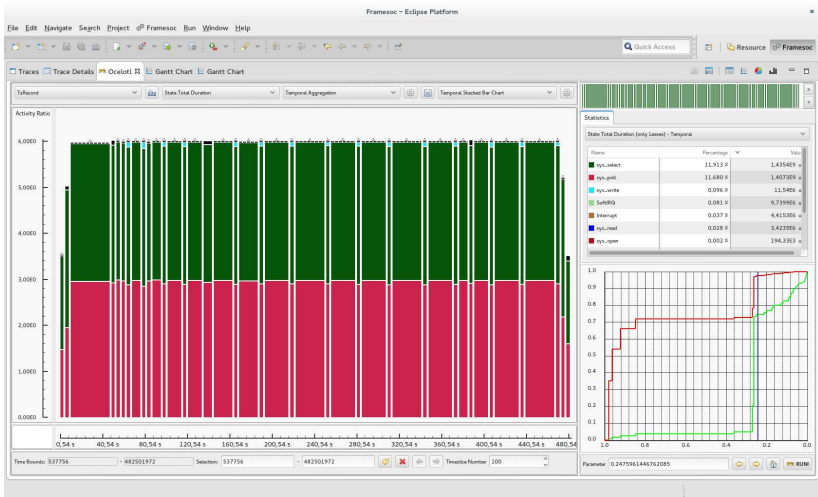
EXPERIMENTS ON EMBEDDED SYSTEMS



$p=0.47$

Multi-level aggregation: Ocelot Application/Demo

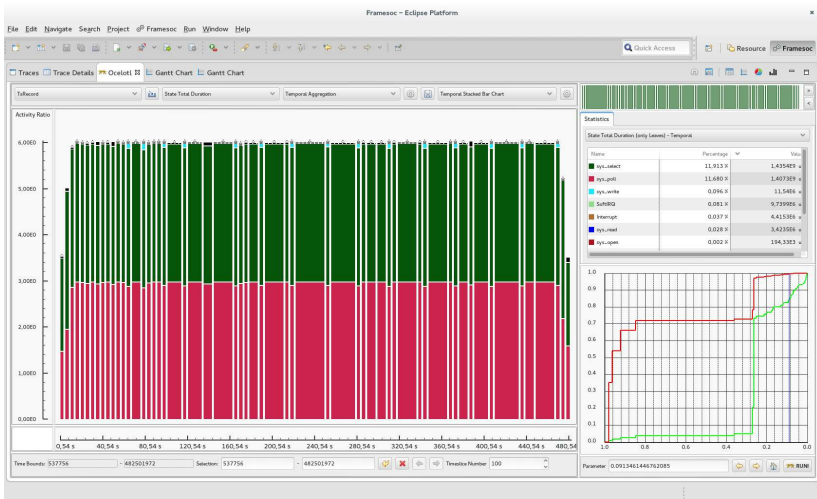
EXPERIMENTS ON EMBEDDED SYSTEMS



$p=0.24$

Multi-level aggregation: Ocelot Application/Demo

EXPERIMENTS ON EMBEDDED SYSTEMS

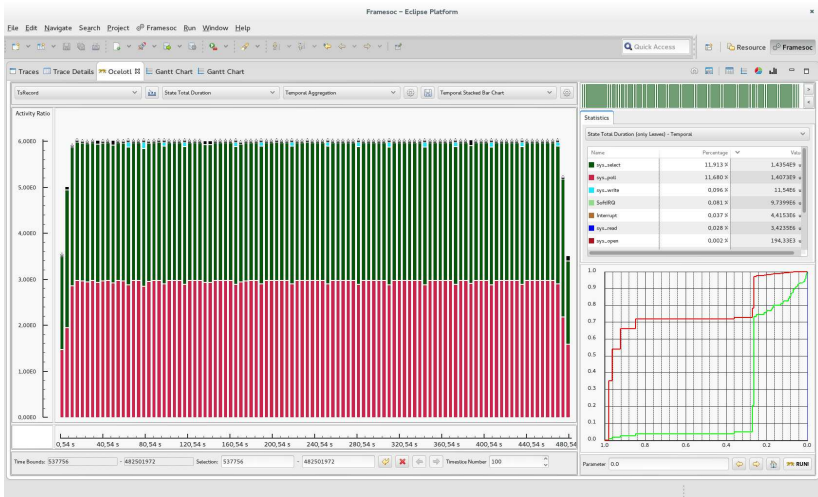


$p=0.09$

Multi-level aggregation: Ocelotl Application/Demo



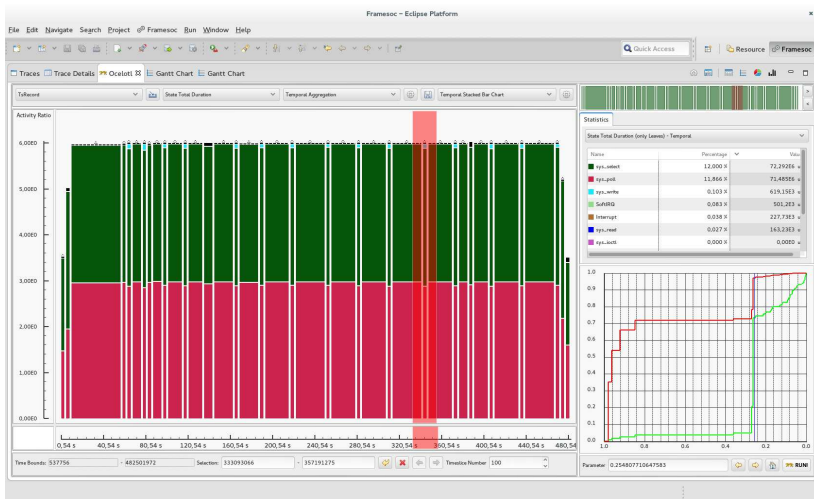
EXPERIMENTS ON EMBEDDED SYSTEMS



$p=0$

Multi-level aggregation: Ocelotl Application/Demo

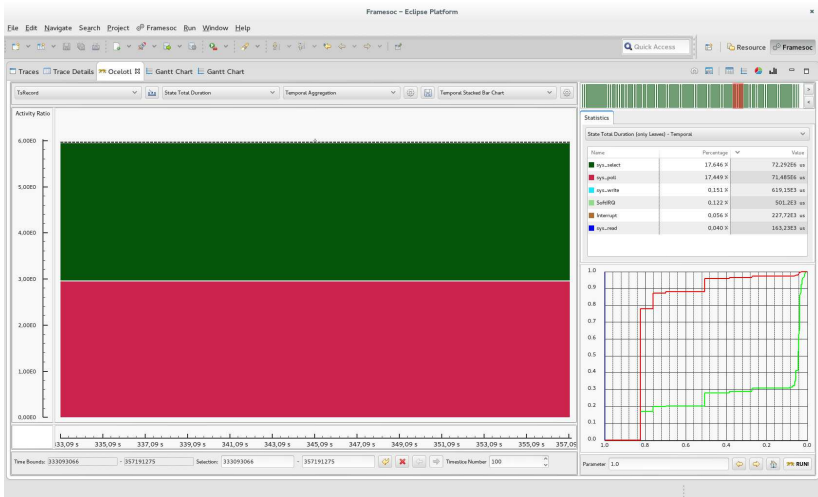
EXPERIMENTS ON EMBEDDED SYSTEMS



$p=0.25$

Multi-level aggregation: Ocelotl Application/Demo

EXPERIMENTS ON EMBEDDED SYSTEMS

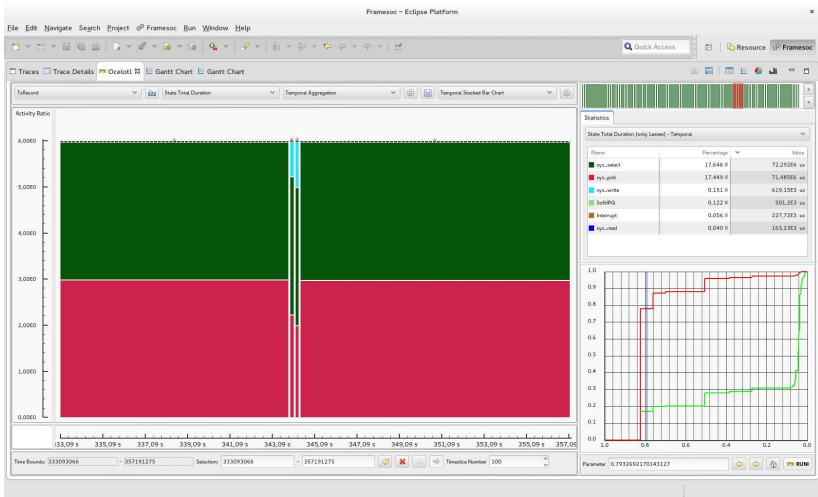


$p=1$

Multi-level aggregation: Ocelot Application/Demo



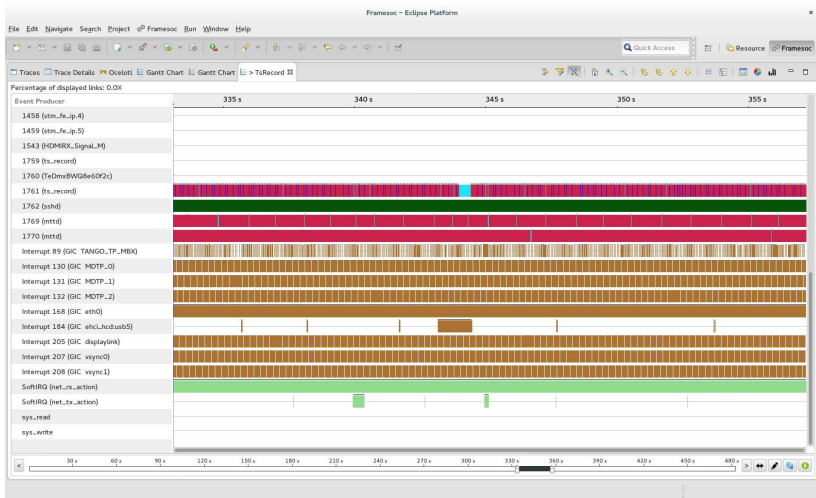
EXPERIMENTS ON EMBEDDED SYSTEMS



$p=0.79$

Multi-level aggregation: Ocelotl Application/Demo

EXPERIMENTS ON EMBEDDED SYSTEMS



All events : `sys_write()` blocked writing on USB disk

Multi-level aggregation: Ocelot! Application/Demo



SPATIOTEMPORAL AGGREGATION

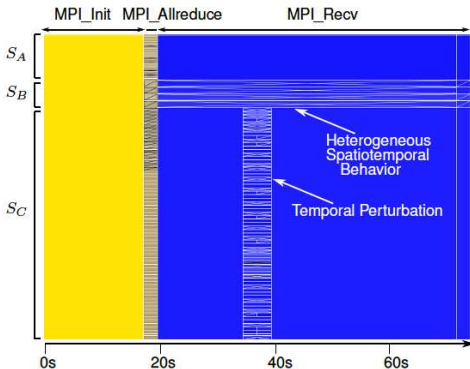


Figure 4. Ocelotl overview of the MPI application LU, class C, 700 processes, executed on the Nancy site of Grid'5000 (S_A : Graphene, S_B : Graphite, S_C : Griffon). We mainly distinguish an initialization sequence (0-20s), followed by the computation phase, where the behavior of the Graphite cluster is heterogeneous in space and time, and there is a perturbation that touches only the execution of the Griffon cluster (34.5s).

SPATIOTEMPORAL AGGREGATION

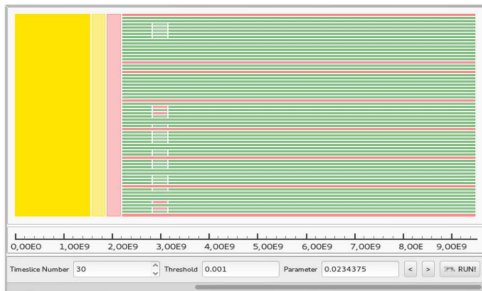


Figure 1. Our analysis tool, Ocelotl, showing an overview of the execution of the NAS-CG application, class C, 64 processes, on the Grid'5000 Rennes site: the trace is partitioned into aggregates that correspond to a locally homogeneous behavior of the application over time and among a set of computing resources. We distinguish a perturbation around 3,00E9, caused by the concurrent execution of applications competing for network access.

HOW TO FIND A NEEDLE IN A HAYSTACK ?

- 1 THE PROBLEM : Extracting macroscopic information from microscopic measures
- 2 ANOMALY : Heterogeneous (unexpected) Macroscopic Behavior
 - SPACE : Localization Heterogeneity
 - TIME : Temporal Heterogeneity
 - SPACE/TIME : Behavior Heterogeneity
- 3 **METHODOLOGY : Information based Aggregation**
 - MACROSCOPIC STATE : Information based approach
 - PARTITION : Algorithms and Complexity
- 4 SYNTHESIS AND OPEN QUESTIONS

BUILDING MACROSCOPIC INFORMATION

The Clustering Approach

- ▶ Similarity of objects (distance function, usually in an euclidian space)
- ▶ Many methods, (k-means, hierarchical,...)
- ▶ Level of clustering (dendograms)

Main difficulties

- ⇒ distance function: semantic of the function
- ⇒ definition of distance between clusters (centroids, center of gravity,...)
- ⇒ semantic of the quality function

Aggregation Approach

Definition of an aggregate

- ▶ set of locations
- ▶ time period
- ▶ probability distribution on the state-space of objects

Restrictions to a set of meaningful subsets

- ▶ contiguous locations
- ▶ time intervals
- ⇒ External information: hierarchy, topology, time contiguity ...

Quality of an aggregation function

- ▶ Information loss
- ▶ Complexity gain
- ⇒ allows comparison, composition, and interpretation (coding)



BUILDING MACROSCOPIC INFORMATION

The Clustering Approach

- ▶ Similarity of objects (distance function, usually in an euclidian space)
- ▶ Many methods, (k-means, hierarchical,...)
- ▶ Level of clustering (dendograms)

Main difficulties

- ⇒ distance function: semantic of the function
- ⇒ definition of distance between clusters (centroids, center of gravity,...)
- ⇒ semantic of the quality function

Aggregation Approach

Definition of an aggregate

- ▶ set of locations
- ▶ time period
- ▶ probability distribution on the state-space of objects

Restrictions to a set of meaningful subsets

- ▶ contiguous locations
- ▶ time intervals
- ⇒ External information: hierarchy, topology, time contiguity ...

Quality of an aggregation function

- ▶ Information loss
- ▶ Complexity gain
- ⇒ allows comparison, composition, and interpretation (coding)



BUILDING MACROSCOPIC INFORMATION

The Clustering Approach

- ▶ Similarity of objects (distance function, usually in an euclidian space)
- ▶ Many methods, (k-means, hierarchical,...)
- ▶ Level of clustering (dendograms)

Main difficulties

- ⇒ distance function: semantic of the function
- ⇒ definition of distance between clusters (centroids, center of gravity,...)
- ⇒ semantic of the quality function

Aggregation Approach

Definition of an aggregate

- ▶ set of locations
- ▶ time period
- ▶ probability distribution on the state-space of objects

Restrictions to a set of meaningful subsets

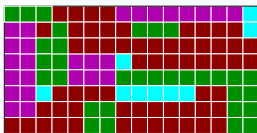
- ▶ contiguous locations
- ▶ time intervals
- ⇒ External information: hierarchy, topology, time contiguity ...

Quality of an aggregation function

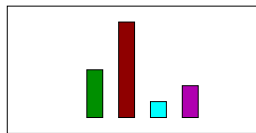
- ▶ Information loss
- ▶ Complexity gain
- ⇒ allows comparison, composition, and interpretation (coding)



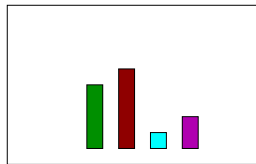
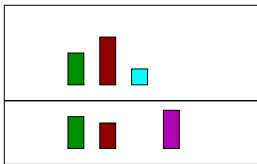
AGGREGATION PROCESS



Microscopic information



Aggregated information

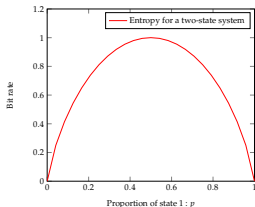


Composition of aggregates

ENTROPY : MEASURE OF HOMOGENEITY/DISORDER

Entropy

$$\begin{aligned}
 H &= - \sum \frac{|s_k|}{|S|} \log_2 \frac{|s_k|}{|S|} \\
 &= \sum p_k \log_2 \frac{1}{p_k}
 \end{aligned}$$



Quantity of information to code the system

Entropy Properties

- ▶ $H \geq 0, H(p) = 0 \Rightarrow$ deterministic system
- ▶ $H(p) \leq \log_2 n, H(p) = \log_2 n \Rightarrow$ uniform system
- ▶ Independence property
- ▶ Conditioning

Entropy Gain

$$G = H_{micro} - H_{macro}$$

- ▶ $G \geq 0$
- ▶ $G = 0$ (no aggregation or deterministic micro-system)
- ▶ maximal if one aggregate
- ▶ Composition property

ENTROPY AND COMPLEXITY

Parametrized Information Criterion

$$pRIC(\mathcal{A}) = p \times Gain(\mathcal{A}) - (1 - p) \times Div(\mathcal{A})$$

Shannon **entropy** measure of complexity

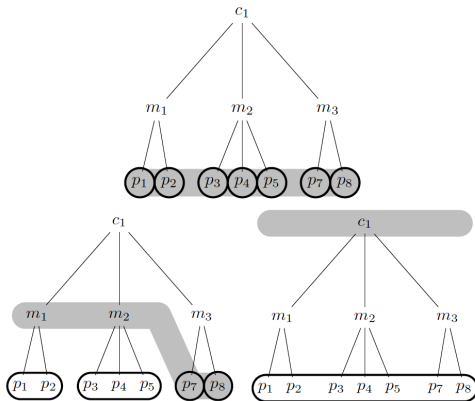
$$Gain(A) = (v(A) \log_2 v(A)) - \sum_{e \in A} (v(e) \log_2 v(e))$$

Kullback-Leibler **divergence** estimates the information loss

$$Div(A) = \sum_{e \in A} v(e) \times \log_2 \left(\frac{v(e)}{v(A)} \times |A| \right)$$

- ▶ The *sum property* of quality measures enables independent computation of aggregates
- ▶ Hierarchical organization allows a recursive evaluation of branches
- ▶ Efficiently implemented through dynamic programming

HIERARCHICAL STRUCTURE



HIERARCHICAL AND TEMPORAL STRUCTURE

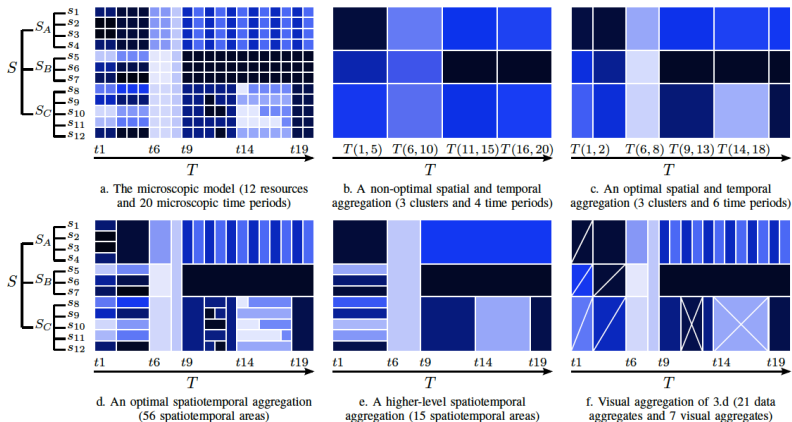


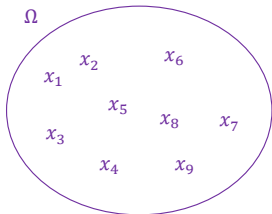
Figure 3. Aggregation and visualization of an artificial trace giving the behavior of 12 resources during 20 microscopic time periods (two possible states)

ALGORITHMS AND COMPLEXITY

The Set Partitioning Problem

Given:

- a set of individuals $\Omega = \{x_1, \dots, x_n\}$

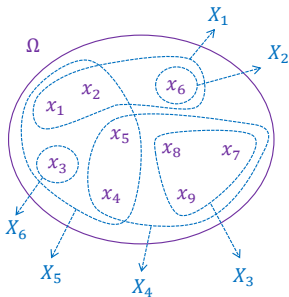


ALGORITHMS AND COMPLEXITY

The Set Partitioning Problem

Given:

- a set of individuals $\Omega = \{x_1, \dots, x_n\}$
- a set of admissible parts $\mathcal{P} = \{X_1, \dots, X_m\} \subset 2^\Omega$

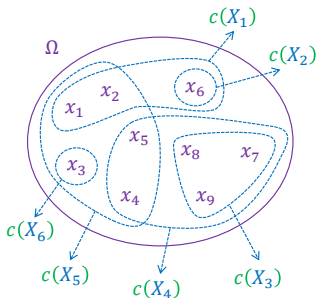


ALGORITHMS AND COMPLEXITY

The Set Partitioning Problem

Given:

- a set of individuals $\Omega = \{x_1, \dots, x_n\}$
- a set of admissible parts $\mathcal{P} = \{X_1, \dots, X_m\} \subset 2^\Omega$
- a cost function $c : \mathcal{P} \rightarrow \mathbb{R}$

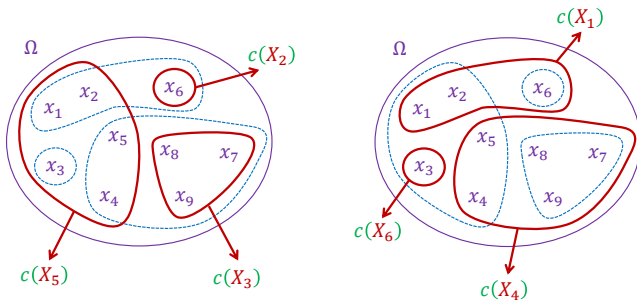


ALGORITHMS AND COMPLEXITY

The Set Partitioning Problem

Given:

- a set of individuals $\Omega = \{x_1, \dots, x_n\}$
- a set of admissible parts $\mathcal{P} = \{X_1, \dots, X_m\} \subset 2^\Omega$
- a cost function $c : \mathcal{P} \rightarrow \mathbb{R}$
- the corresponding set of admissible partitions $\mathfrak{P} = \{\mathcal{X} \subset \mathcal{P} \text{ such that } \mathcal{X} \text{ is a partition of } \Omega\}$

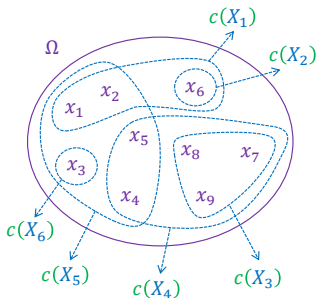


ALGORITHMS AND COMPLEXITY

The Set Partitioning Problem

Given:

- a set of individuals $\Omega = \{x_1, \dots, x_n\}$
- a set of admissible parts $\mathcal{P} = \{X_1, \dots, X_m\} \subset 2^\Omega$
- a cost function $c : \mathcal{P} \rightarrow \mathbb{R}$
- the corresponding set of admissible partitions $\mathfrak{P} = \{\mathcal{X} \subset \mathcal{P} \text{ such that } \mathcal{X} \text{ is a partition of } \Omega\}$



Problem: Find an admissible partition that minimizes the cost function:

$$\mathcal{X}^* = \arg \min_{\mathcal{X} \in \mathfrak{P}} \left(\sum_{X \in \mathcal{X}} c(X) \right)$$

→ NP-complete!

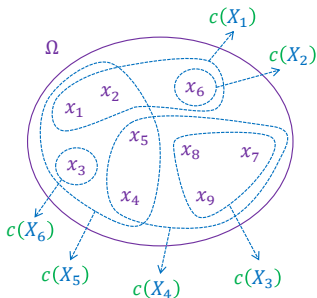
ALGORITHMS AND COMPLEXITY

The Set Partitioning Problem

Given:

- a set of individuals $\Omega = \{x_1, \dots, x_n\}$
- a set of admissible parts $\mathcal{P} = \{X_1, \dots, X_m\} \subset 2^\Omega$
- a cost function $c : \mathcal{P} \rightarrow \mathbb{R}$
- the corresponding set of admissible partitions $\mathfrak{P} = \{\mathcal{X} \subset \mathcal{P} \text{ such that } \mathcal{X} \text{ is a partition of } \Omega\}$

Additional assumptions



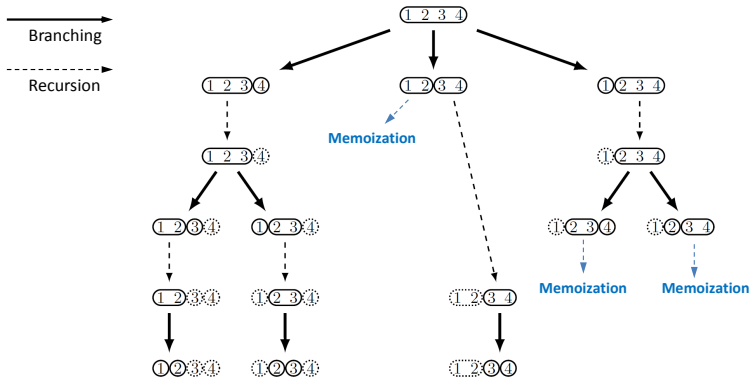
Problem: Find an admissible partition that minimizes the cost function:

$$\mathcal{X}^* = \arg \min_{\mathcal{X} \in \mathfrak{P}} \left(\sum_{X \in \mathcal{X}} c(X) \right)$$

→ NP-complete!

ALGORITHMS AND COMPLEXITY

Memoization



HOW TO FIND A NEEDLE IN A HAYSTACK ?

- 1 **THE PROBLEM** : Extracting macroscopic information from microscopic measures
- 2 **ANOMALY** : Heterogeneous (unexpected) Macroscopic Behavior
 - SPACE : Localization Heterogeneity
 - TIME : Temporal Heterogeneity
 - SPACE/TIME : Behavior Heterogeneity
- 3 **METHODOLOGY** : Information based Aggregation
 - MACROSCOPIC STATE : Information based approach
 - PARTITION : Algorithms and Complexity
- 4 **SYNTHESIS AND OPEN QUESTIONS**

SYNTHESIS AND OPEN QUESTIONS

Proposition: Multi-scale aggregation driven by information content

- ▶ Automatic computation : information contained in a representation
- ▶ Optimal representation (compromise complexity/information loss) according a hierarchy structure
- ▶ **Polynomial** time computation of the optimal tree : $\mathcal{O}(n)$, sequence : $\mathcal{O}(n^2)$, spatio-temporal $\mathcal{O}(n^3)$
- ▶ Scalable now to a million

Open-source software – **Viva** and **Ocelotl**

- ▶ <https://github.com/schnorr/viva>
- ▶ <https://soctrace-inria.github.io/ocelotl/>
- ▶ Currently being packaged to Debian

Issues yet to be addressed

- ▶ Generalization of the approach/algorithm to : networks of processes (graph based aggregation), geometrically located processes, sets of causally-related events (causal aggregation), semantic aggregation (superstates).
- ▶ Combination with other complexity measures : Minimal Description Length, theoretical framework (Kolmogorov Complexity)
- ▶ Application to other scientific domains: embedded systems, multi-agent systems, geography, social sciences...



SYNTHESIS AND OPEN QUESTIONS

Proposition: Multi-scale aggregation driven by information content

- ▶ Automatic computation : information contained in a representation
- ▶ Optimal representation (compromise complexity/information loss) according a hierarchy structure
- ▶ **Polynomial** time computation of the optimal tree : $\mathcal{O}(n)$, sequence : $\mathcal{O}(n^2)$, spatio-temporal $\mathcal{O}(n^3)$
- ▶ Scalable now to a million

Open-source software – **Viva** and **Ocelotl**

- ▶ <https://github.com/schnorr/viva>
- ▶ <https://soctrace-inria.github.io/ocelotl/>
- ▶ Currently being packaged to Debian

Issues yet to be addressed

- ▶ Generalization of the approach/algorithm to : networks of processes (graph based aggregation), geometrically located processes, sets of causally-related events (causal aggregation), semantic aggregation (superstates).
- ▶ Combination with other complexity measures : Minimal Description Length, theoretical framework (Kolmogorov Complexity)
- ▶ Application to other scientific domains: embedded systems, multi-agent systems, geography, social sciences...



SYNTHESIS AND OPEN QUESTIONS

Proposition: Multi-scale aggregation driven by information content

- ▶ Automatic computation : information contained in a representation
- ▶ Optimal representation (compromise complexity/information loss) according a hierarchy structure
- ▶ **Polynomial** time computation of the optimal tree : $\mathcal{O}(n)$, sequence : $\mathcal{O}(n^2)$, spatio-temporal $\mathcal{O}(n^3)$
- ▶ Scalable now to a million

Open-source software – **Viva** and **Ocelotl**

- ▶ <https://github.com/schnorr/viva>
- ▶ <https://soctrace-inria.github.io/ocelotl/>
- ▶ Currently being packaged to Debian

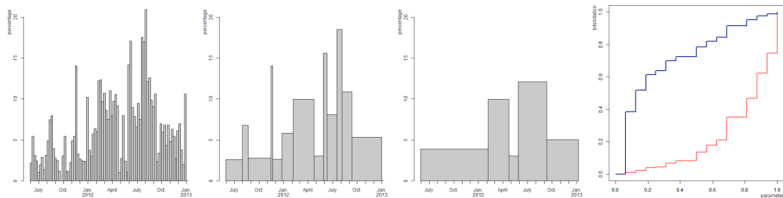
Issues yet to be addressed

- ▶ Generalization of the approach/algorithm to : networks of processes (graph based aggregation), geometrically located processes, sets of causally-related events (causal aggregation), semantic aggregation (superstates).
- ▶ Combination with other complexity measures : Minimal Description Length, theoretical framework (Kolmogorov Complexity)
- ▶ Application to other scientific domains: embedded systems, multi-agent systems, geography, social sciences...



APPLICATION : INTERNATIONAL RELATIONS THROUGH NEWSPAPERS RSS

Figure 1: Time aggregation of the probability that Syria appears in the RSS flows "Times of India - International"



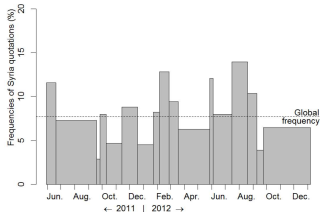
APPLICATION : INTERNATIONAL RELATIONS THROUGH NEWSPAPERS RSS

Figure 3: Weekly frequencies of Syria quotations in four RSS flows

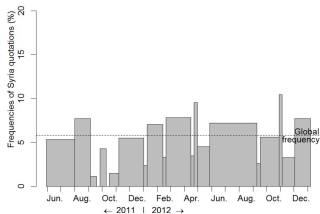


APPLICATION : INTERNATIONAL RELATIONS THROUGH NEWSPAPERS RSS

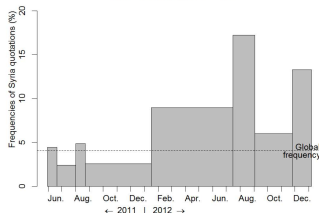
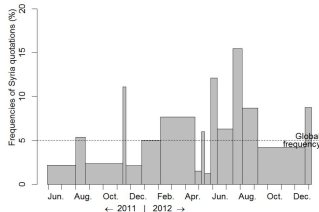
Figure 4: Optimal time aggregation of the frequencies of Syria quotations in four RSS flows
Le Monde
The Times of India



The Financial Times



The Washington Post



SPATIAL (HIERARCHICAL) AGGREGATION

IJCAI 2013 VIDEO COMPETITION

August 3-9, 2013, Beijing

Multi-resolution Representations of Media Information

R. Lamarche-Perrin

Y. Demazeau

J.-M. Vincent



LABORATOIRE D'INFORMATIQUE DE GRENOBLE

RSS - flows \Rightarrow international relations

SOME REFERENCES

- ▶ **A Spatiotemporal Aggregation Technique for Overviewing Large Traces** Damien Dosimont et al. Cluster 2014, Madrid
- ▶ **Building Optimal Macroscopic Representations of Complex Multi-agent Systems** Robin Lamarche-Perrin, et al. T. Computational Collective Intelligence 15: 1-27 (2014)
- ▶ **Emergence and Macroscopic Observation of Large-scale Multi-agent Systems** Robin Lamarche-Perrin et al. Bracis 2014, São Carlos
- ▶ **Building Optimal Macroscopic Representations of Complex Multi-agent Systems. Application to the Spatial and Temporal Analysis of International Relations through News Aggregation.** Robin Lamarche-Perrin et al. In Trans. on Comput. Collective Intelligence, 2014
- ▶ **Agrégation de traces pour la visualisation de grands systèmes distribués** Robin Lamarche-Perrin et al. TSI, 2014
- ▶ **Evaluating Trace Aggregation for Performance Visualization of Large Distributed Systems** Robin Lamarche-Perrin et al. ISPASS'14, Monterey 2014.
- ▶ **The Best-partitions Problem: How to Build Meaningful Aggregations** IAT, 2013, Atlanta.
- ▶ **Identification of International Media Events by Spatial and Temporal Aggregation of RSS Flows of Newspapers. Application to the Case of the Syrian Civil War between May 2011 and December 2012.** Timothe Giraud et al. ECTQG, 2013, Dourdan.
- ▶ **How to Build the Best Macroscopic Description of your Multi-agent System?** Robin Lamarche-Perrin et al. PAAMS'13, Salamanca.
- ▶ **Macroscopic Analysis of Large-scale systems. Epistemic Emergence and Spatiotemporal Aggregation** Robin Lamarche-Perrin Ph.D, University of Grenoble 2014
- ▶ **Agrégation spatiotemporelle pour la visualisation de traces d'exécution** Ph.D, University of Grenoble Damien Dosimont 2015

STOCHASTIC PROCESSES AND ENTROPY

Conditioning

Entropy of a couple of r.v. (X, Y)

$$H(X, Y) = H(Y) + H(X|Y)$$

X and Y are independent iff

$$H(X, Y) = H(X) + H(Y).$$

Entropy rate of a stochastic process

$$H(X) = \lim_n \frac{1}{n} H(X_1, \dots, X_n)$$

Conditional entropy rate of a stochastic process

$$H'(X) = \lim_n \frac{1}{n} H(X_n | X_{n-1}, \dots, X_1)$$

Remark (stationary):

$H(X_n | X_{n-1}, \dots, X_1)$ is nonincreasing

Theorem

For a stationary process

$$H(X) = H'(X)$$

i.i.d. process : $\{X_n\}$

$$H(X) = H'(X) = H(X_1)$$

Homogeneous stationary Markov chain

$$H(X) = H(X_2 | X_1) = \sum_i \pi_i \sum_j p_{i,j} \log \frac{1}{p_{i,j}}$$

References

- ▶ A Mathematical Theory of Communication
C.E. Shannon The Bell System Technical Journal. 1948
- ▶ Elements of Information Theory *T. M. Cover, J. A. Thomas* John Wiley & Sons, Inc. 2006
- ▶ Entropy and Information Theory *R. Gray*



STOCHASTIC PROCESSES AND ENTROPY

Conditioning

Entropy of a couple of r.v. (X, Y)

$$H(X, Y) = H(Y) + H(X|Y)$$

X and Y are independent iff

$$H(X, Y) = H(X) + H(Y).$$

Entropy rate of a stochastic process

$$H(X) = \lim_n \frac{1}{n} H(X_1, \dots, X_n)$$

Conditional entropy rate of a stochastic process

$$H'(X) = \lim_n \frac{1}{n} H(X_n | X_{n-1}, \dots, X_1)$$

Remark (stationary):

$H(X_n | X_{n-1}, \dots, X_1)$ is nonincreasing

Theorem

For a stationary process

$$H(X) = H'(X)$$

i.i.d. process : $\{X_n\}$

$$H(X) = H'(X) = H(X_1)$$

Homogeneous stationary Markov chain

$$H(X) = H(X_2 | X_1) = \sum_i \pi_i \sum_j p_{i,j} \log \frac{1}{p_{i,j}}$$

References

- ▶ A Mathematical Theory of Communication
C.E. Shannon The Bell System Technical Journal. 1948
- ▶ Elements of Information Theory T. M. Cover, J. A. Thomas John Wiley & Sons, Inc. 2006
- ▶ Entropy and Information Theory R. Gray

