QUDOS 2016
Saarbrücken, Germany

# A Tool for Verification of Big-Data Applications

## *Jul 21th, 2016*

M.M. Bersani, F. Marconi, M.G. Rossi
Politecnico di Milano
Milan, Italy

Madalina Erascu
Institute e-Austria Timisoara &
Western University of Timisoara
Timisoara, Romania

# Our work at a glance

o Approach and tool for the automated verification of topology-based data-intensive applications.

- Based (so far) on temporal logic model

- Performs automated transformation from high level application description to formal model

- Enables verification of safety properties

# Roadmap

- *Context*
  - *Quality assurance in DIA*
- *Research Design*
  - *Research question*
  - *Our approach*
- *Conclusions*
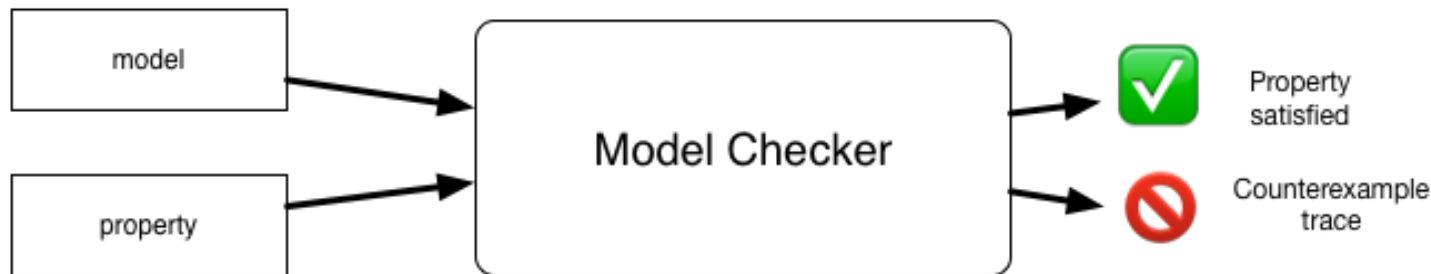  - *Contributions*
  - *Future works*

*Quality Analysis and Verification for data-intensive applications*
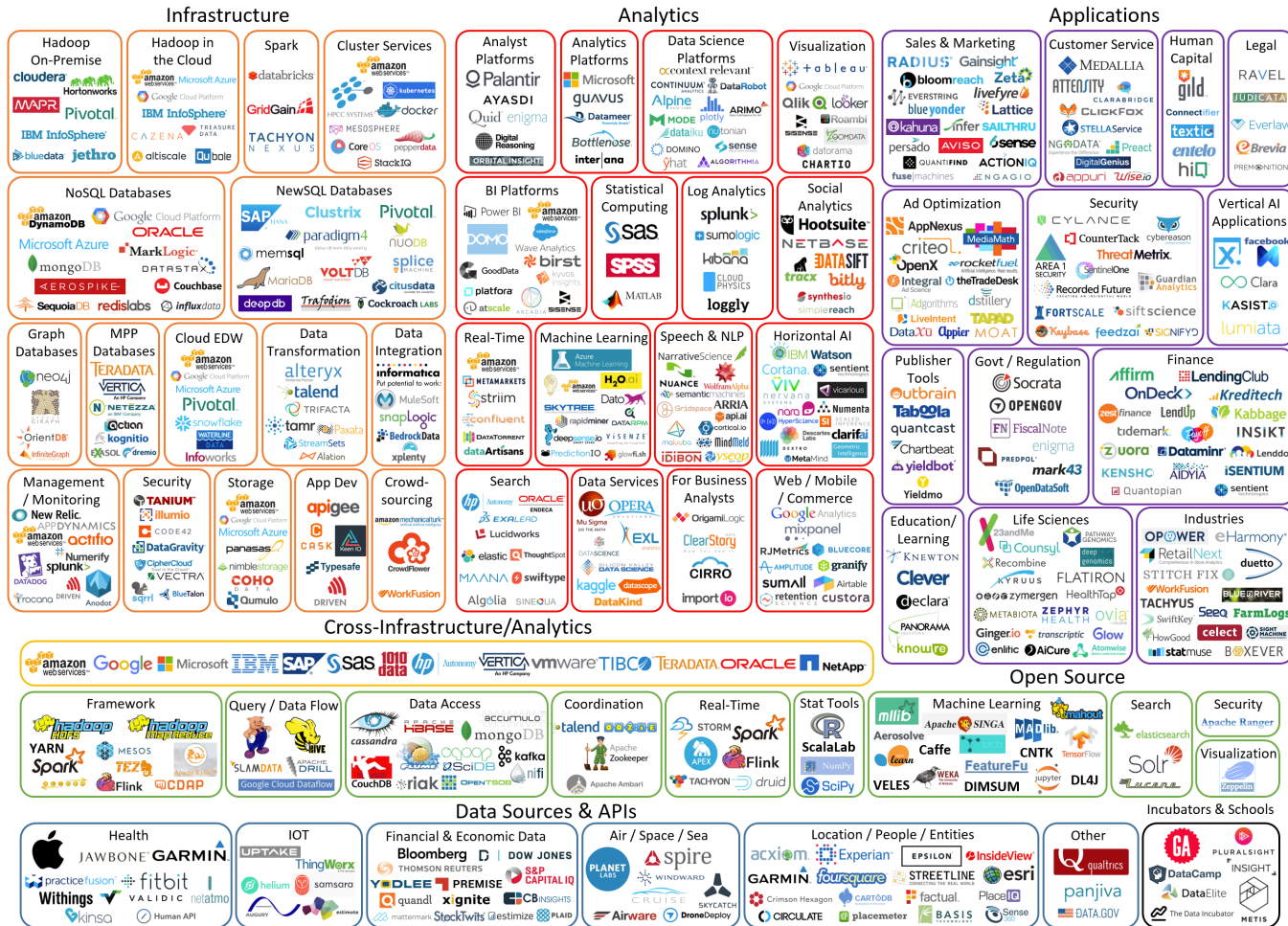
# CONTEXT

# Formal Verification

o Given a Model M and a Property specification P, verification checks whether P holds in M.

o M and P can be expressed in many different ways
  ▪ various kinds of automata (operational models)
  ▪ various kinds of logics (descriptive models)

# Data-Intensive Applications (DIA)



Big Data Landscape 2016 (Version 3.0)

# DICE Project

o Horizon 2020 Research & Innovation Action (RIA)

- Quality-Aware Development for Big Data applications

- Feb 2015 - Jan 2018, 4M Euros budget

- 9 partners (Academia & SMEs), 7 EU countries
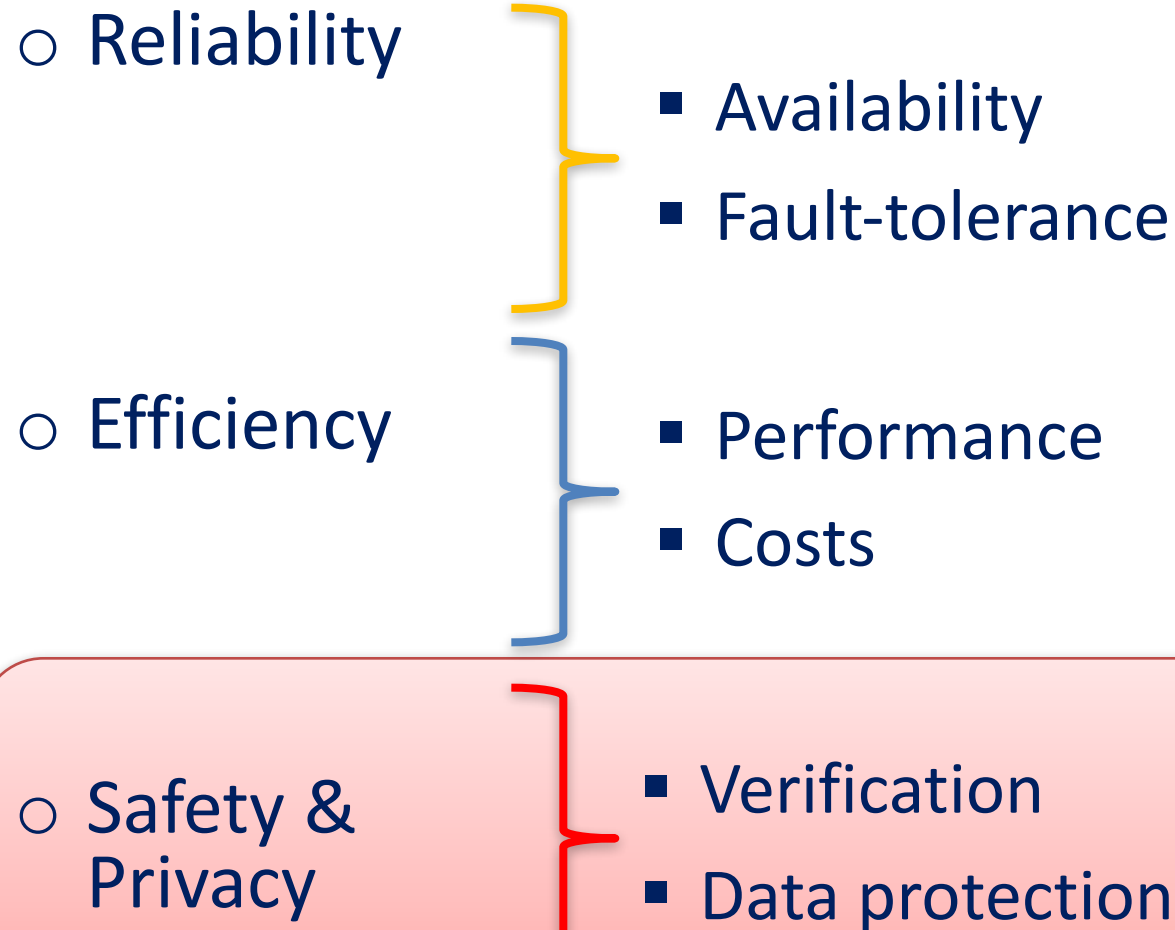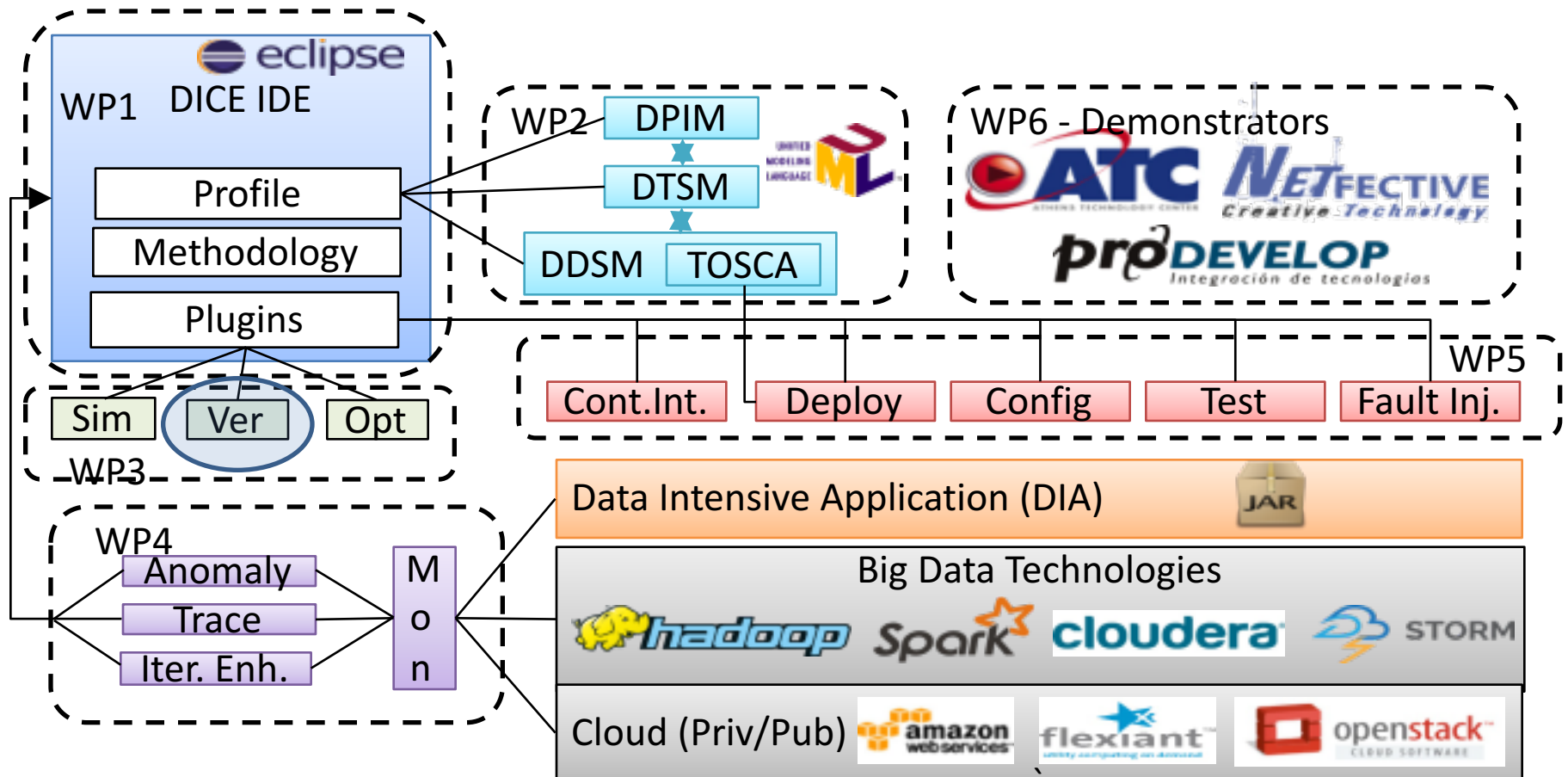
Imperial College London

XLAB NOT IDLE

@ IeAT

Universidad Zaragoza 1542

NETFECTIVE Creative Technology

ATC ATHENS TECHNOLOGY CENTER

POLITECNICO DI MILANO

flexiant TM your cloud simplified

proDEVELOP Integración de tecnologías

# Quality Dimensions in DICE

o Reliability
- Availability
- Fault-tolerance

o Efficiency
- Performance
- Costs

o Safety & Privacy
- Verification
- Data protection
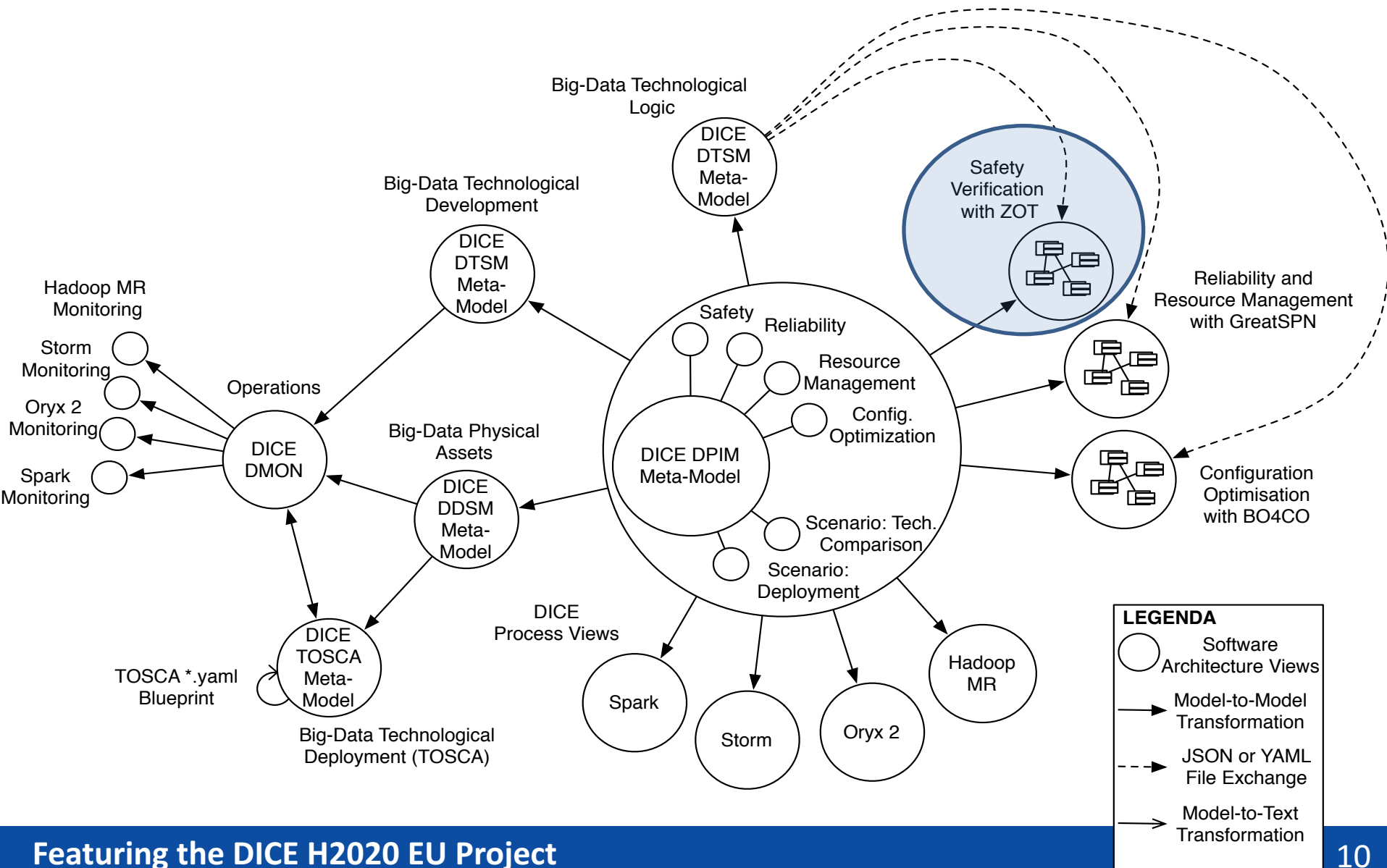
# Our positioning in DICE framework (1)

# Our positioning in DICE framework (2)

*Quality Analysis and Verification for data-intensive applications*

# RESEARCH DESIGN

# Research question

"How can we verify safety properties
of a data-intensive application?"

# State of the art

o Formal verification of distributed systems is a major research area in software engineering

o Few works trying to address formal verification in the context of DIA

- Main focus on verifying *application-independent* properties related to specific frameworks

  - Reliability and load balancing of MapReduce
  - Validity of messaging flow in MapReduce

- no modeling and verification of *application-dependent* properties

o Verification tools have been used as verification engines to build formal verification techniques for UML models

- Few of them deal with real-time constraints.
- Mainly focused on functional requirements.

# Our Approach

- Focus on a specific set of technologies
  - <u>Topology-based streaming applications</u> → **Apache Storm**
- Identify safety issues
- Devise a formal model
  - Having an appropriate level of abstraction
  - Allowing to capture meaningful system behavior and properties
  - Using a formalism that enables automatic verification
- Define a tool-supported mechanism for formal verification
  - Starting from high level application description (annotated UML)

# Apache Storm

o Open Source Distributed Stream Processing System

o Analytics, Log Event processing, etc..

o Reliability, at-least-one semantics

o Wide adoption in production

o Main concepts
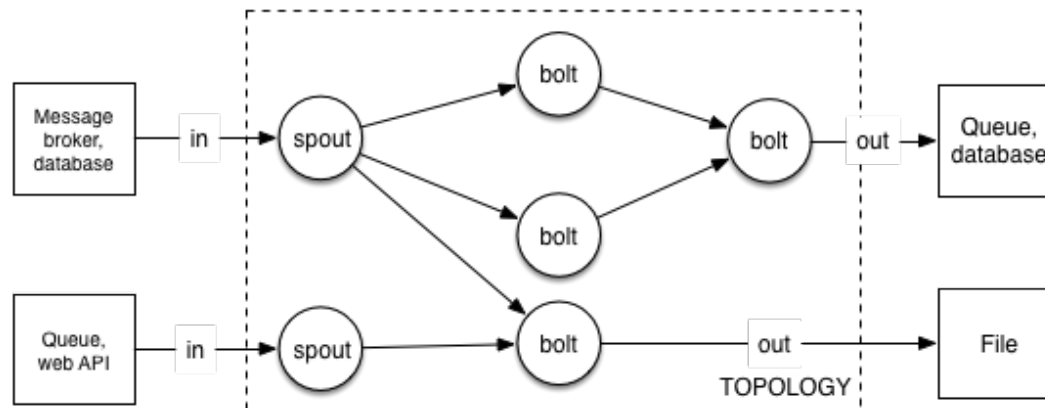
- Streams
- Topologies

# Storm Applications

o Applications defined by means of **Topologies,** graphs of computations composed of:

- **Spouts**
  - Sources of data streams (tuples)
- **Bolts**
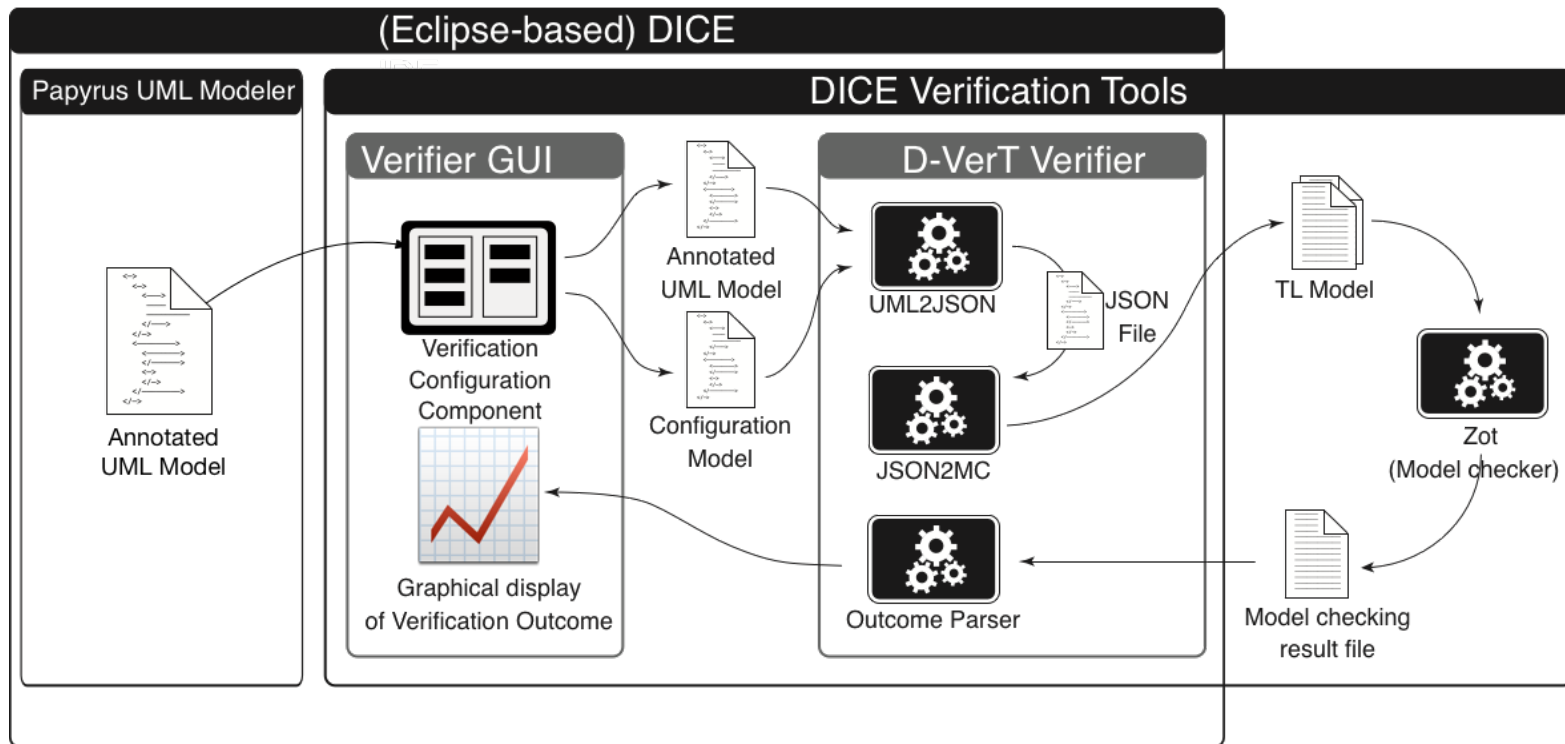  - Calculate, Filter, Aggregate, Join, Talk to databases

# Safety Issues

o Important requirements for streaming applications

- **Latency**
- Throughput

o Critical points

- incorrect design of timing constraints
- node failures

o might cause

- latency in processing tuples
- monotonic growth of the size of used memory (queues).
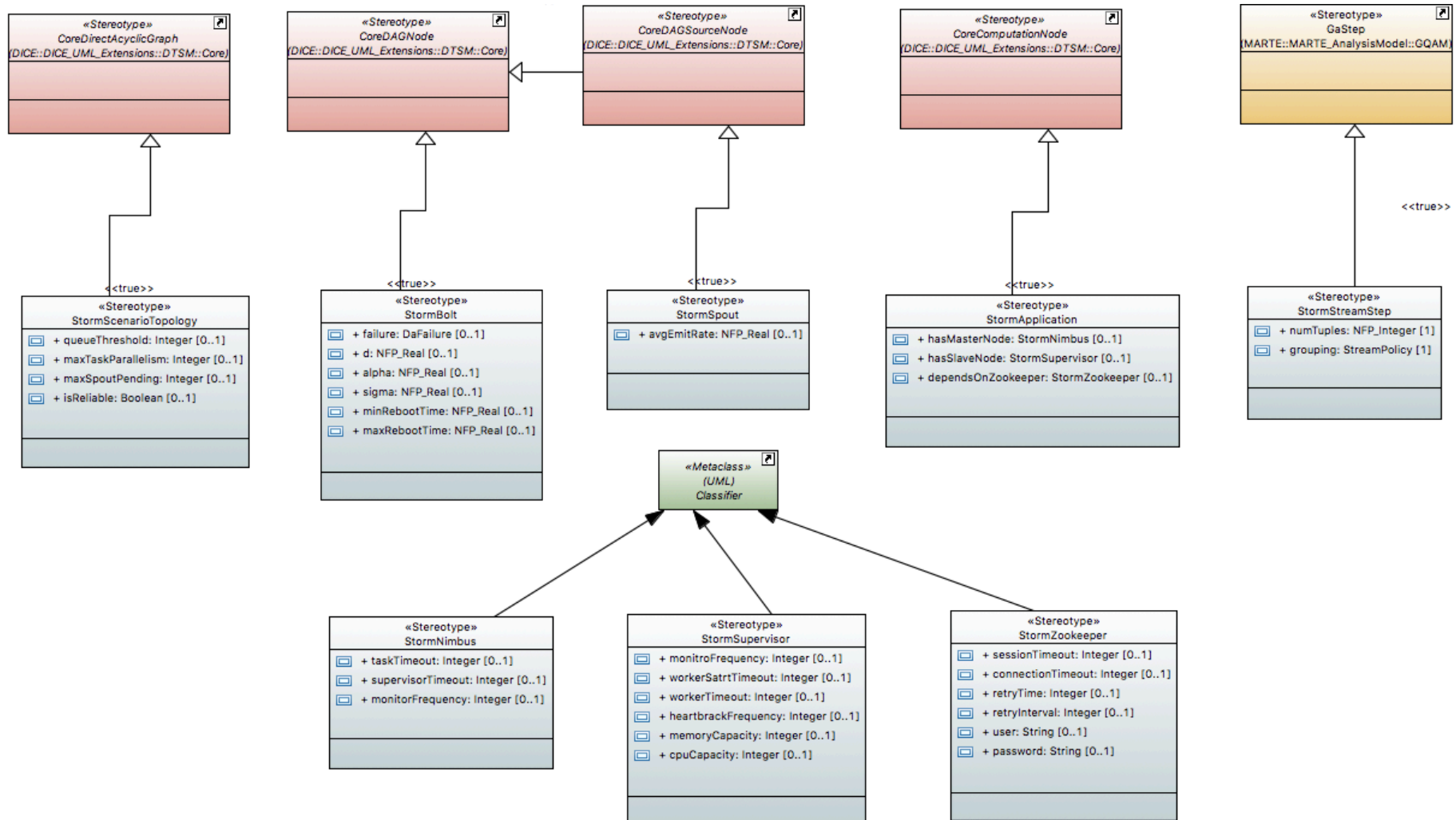
# DICE Verification Tool

o We want to

- Verify whether a topology reaches an **unwanted configuration**

  - e.g., where bolts are not able to process incoming tuples on time

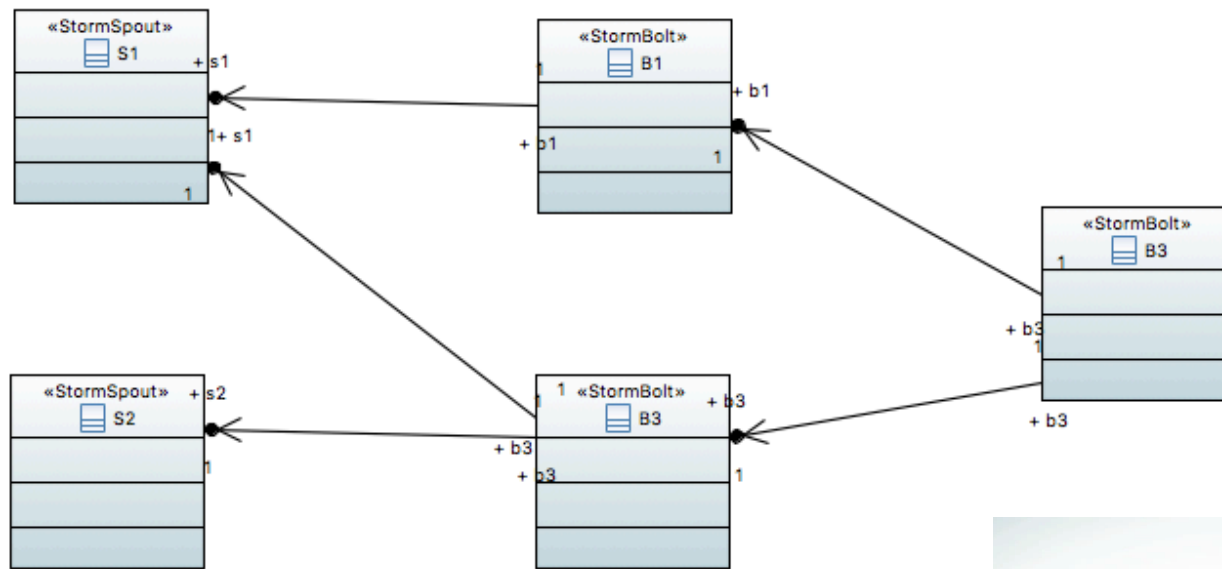- Let the user specify the topology by means of high level models (UML)

# D-VerT - DICE Verification Tool

# DTSM2Json module

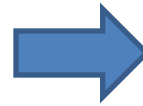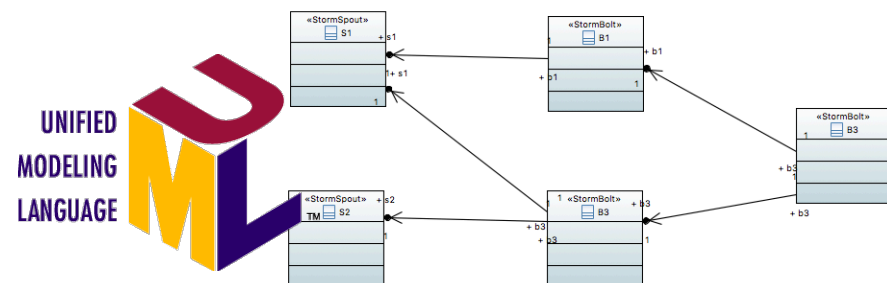o Relies on Eclipse **UML2** Java library

o "Navigates" DTSM class diagram and extract topology structure and information

o Gathers verification option from Eclipse launch configuration

o Maps topology components to Java objects

o Directly converts Java objects to JSON object via **gson** library
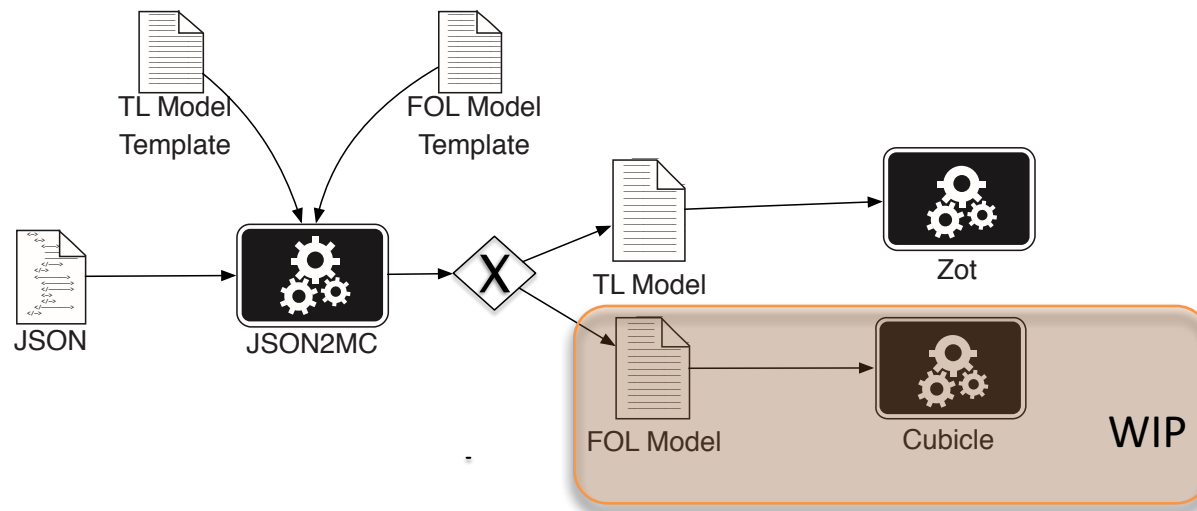
# Json2MC - Module

o Python component based on Jinja2 templating engine

o Generates Formal Model based on the content of JSON file and on the selected template (TL or FOL).
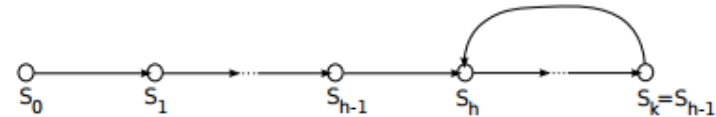
# Verification Approaches

o Bounded Satisfiability Checking (BSC)
- Input:
  - Temporal logic formula (Model)
  - Negated Property over time
- Outcome:
  - SAT → counterexample trace
  - UNSAT → Property holds for the considered time bound
- We use **Zot** verification tool (https://github.com/fm-polimi/zot)

o Reachability Checking (WIP)
- Model defined by FOL Array based system
  - Set of **initial states** and **transitions**
  - Formula defining **undesired states (Negated property)**
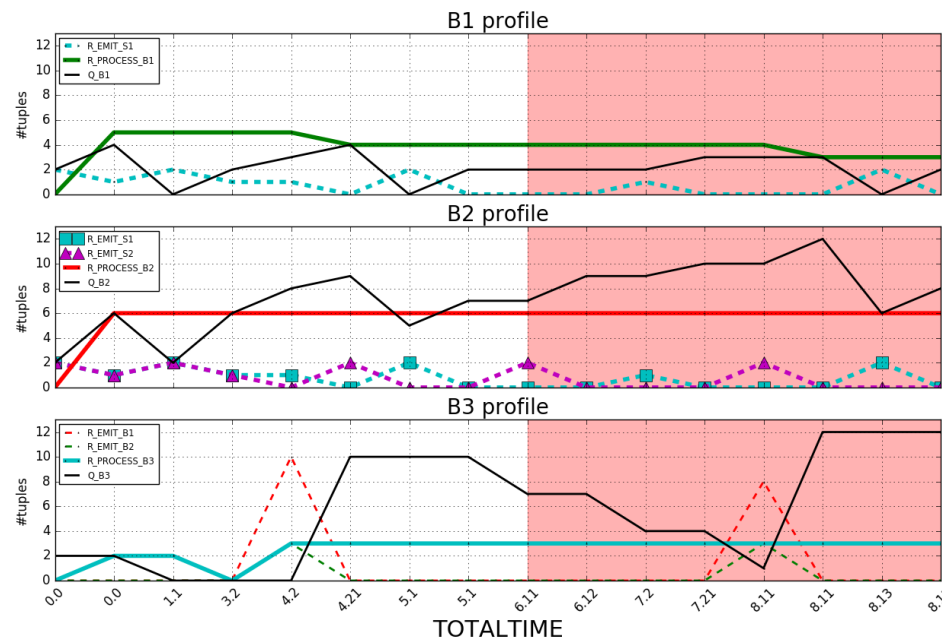- Outcome:
  - UNSAFE → Trace showing that undesired state are reachable from initial states
  - SAFE → No undesired state can be reache from initial states

# D-VerT – Output trace

o When at least one queue grows with an unbounded trend
  - an infinite ultimately periodic model is found
  - **Output Parser** provides graphical counterexample trace

# CONCLUSIONS

# Contributions

o We enabled automatic verification on *topology-based* streaming applications by

- Defining a formal model based on temporal logic

- defining automatic mechanisms for translating to the formal model from a high level description.

- extending Zot Verification tool to support the formalism and carry out BSC on it

# Preliminary results

o Validation through open source and industrial use cases

- Meaningful qualitative results in identifying critical points in topology design

- Execution time strongly depends on the size of the topology and on the configurations of single components

| Topology | Bolts | Time | Max Memory | Outcome | Spurious |
|---|---|---|---|---|---|
| simple-DIA-cfg-1 | 3 | 60s | 104MB | SAT | no |
| simple-DIA-cfg-2 | 3 | 1058s | 150MB | UNSAT | N/A |
| focused-crawler-complete | 8 | 2664s | 448MB | SAT | no |
| focused-crawler-reduced-cfg-1 | 4 | 95s | 142MB | SAT | no |
| focused-crawler-reduced-cfg-2 | 4 | 253s | 195MB | SAT | no |
| focused-crawler-reduced-cfg-3 | 4 | 327s | 215MB | SAT | no |
| focused-crawler-reduced-cfg-4 | 4 | 333s | 206MB | SAT | no |
| focused-crawler-reduced-cfg-5 | 4 | 3184s | 317MB | SAT | yes |
| focused-crawler-reduced-cfg-6 | 4 | 1060s | 229MB | SAT | yes |

http://dice-project.github.io/DICE-Verification/

# Ongoing and Future works

o Identification and verification of further properties
  - Privacy and Security
o Tool improvements
o Modeling different technologies (Spark, CEP, Tez)
o Developing FOL model
o New theoretical results on the correctness and completeness of the formal analysis

# Questions?

Thank you!