



MINING KNOWLEDGE GRAPHS FROM LOOSELY STRUCTURED PROCESSES

A USE CASE FROM EMAILING SYSTEMS

WALID GAALOUL

ICSBT 2022 KEYNOTE

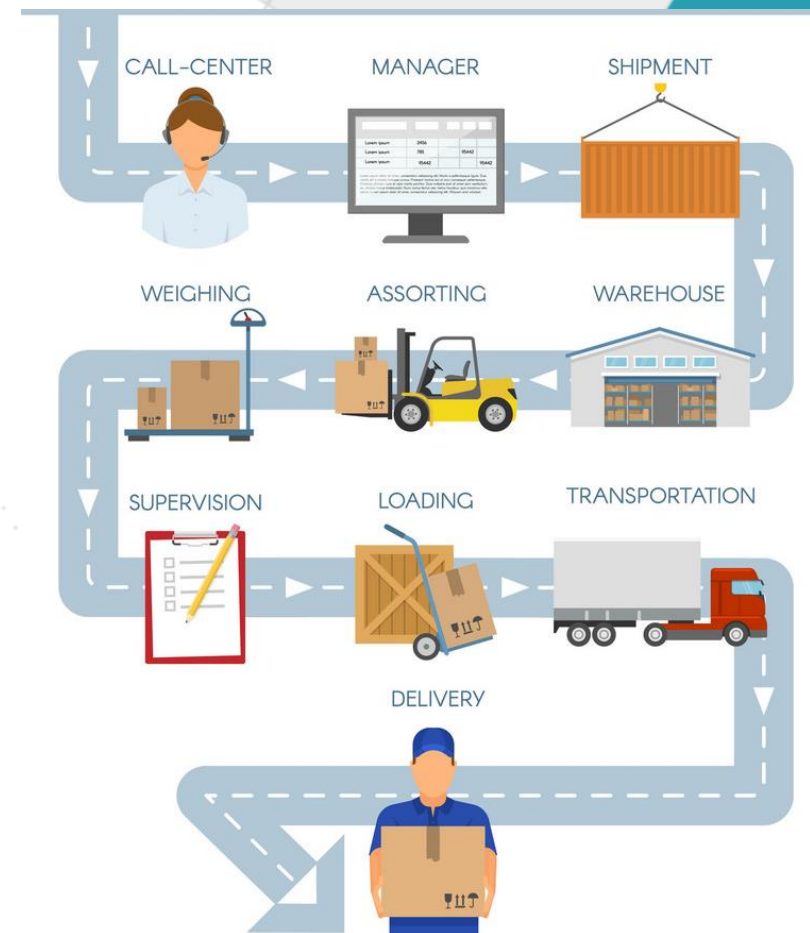
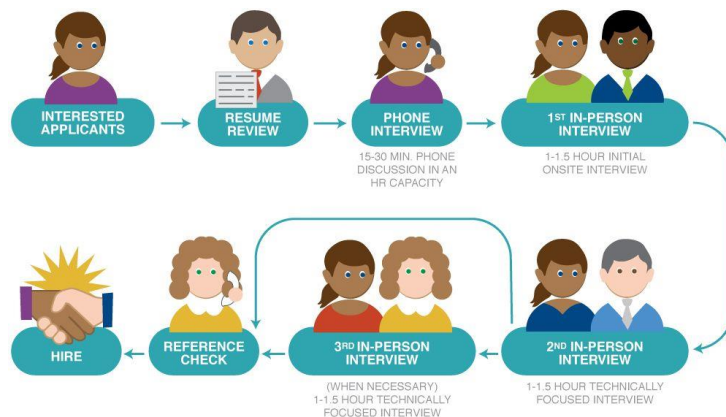


Agenda

- Introduction to Processes & Process Analytics
- Spectrum of processes
 - From structured to unstructured processes
- Conceptual models to enable Cognitive Process Analytics
 - Process Knowledge Graph
 - Conversational Assistants
- Domain specific models for Emailing Systems
 - Discover Process Knowledge Graph from emails
 - Conversational Assistant to query Process Knowledge Graph

What is a process

- A **process** refers to **how work** is done
 - A set of **inter-related activities**
 - involving a number of **actors** and **data**
 - triggered by a need and leading to an **outcome**



WHAT IS A PROCESS

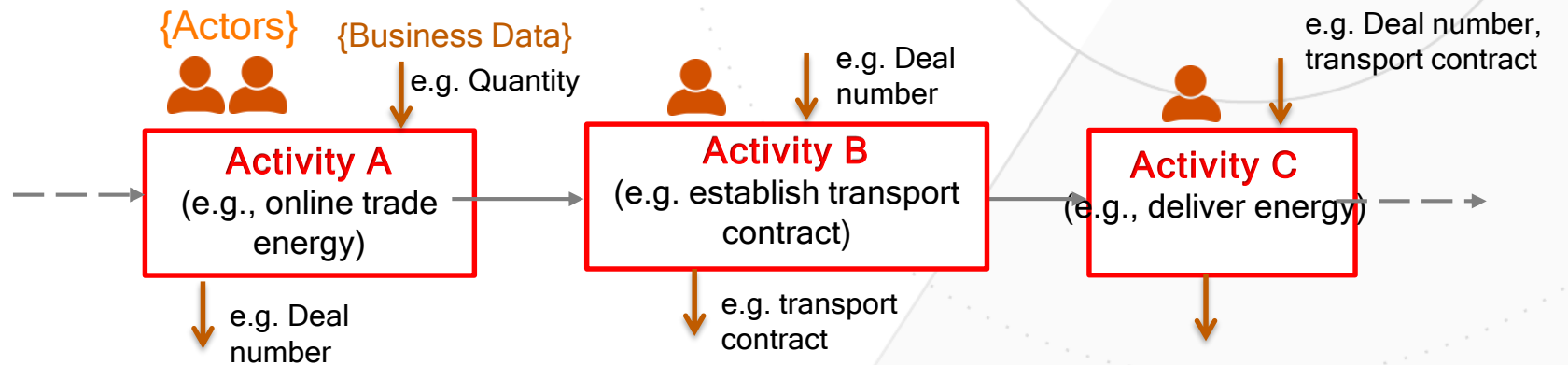
Business Process (BP)

Functional Perspective

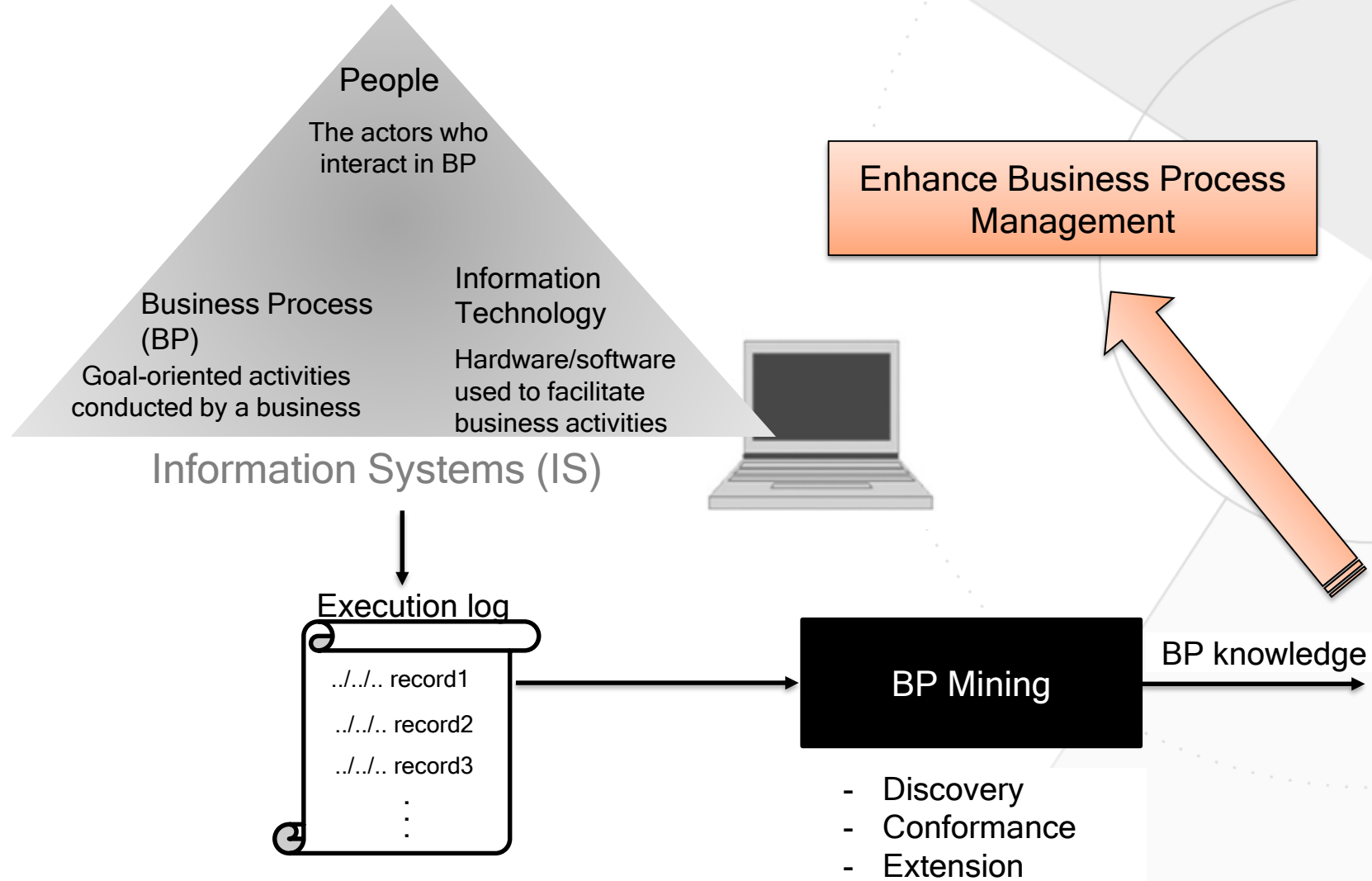
Behavioral Perspective

Organizational Perspective

Data Perspective



PROCESS ANALYTICS LANDSCAPE



Process Execution data

- **Who (resource) did what (activity) and when (timestamp)**
- **What data, objects, artifacts are consumed/produced?**
- 2 standards: **XES**
 - Predefined process instance notion

patient	activity	timestamp	doctor	age	cost
5781	make X-ray	23-1-2014@10.30	Dr. Jones	45	70.00
5541	blood test	23-1-2014@10.18	Dr. Scott	61	40.00
5833	blood test	23-1-2014@10.27	Dr. Scott	24	40.00
5781	blood test	23-1-2014@10.49	Dr. Scott	45	40.00
5781	CT scan	23-1-2014@11.10	Dr. Fox	45	1200.00
5833	surgery	23-1-2014@12.34	Dr. Scott	24	2300.00
5781	handle payment	23-1-2014@12.41	Carol Hope	45	0.00
5541	radiation therapy	23-1-2014@13.57	Dr. Jones	61	140.00
5541	radiation therapy	23-1-2014@13.08	Dr. Jones	61	140.00
..

case id

activity name

timestamp

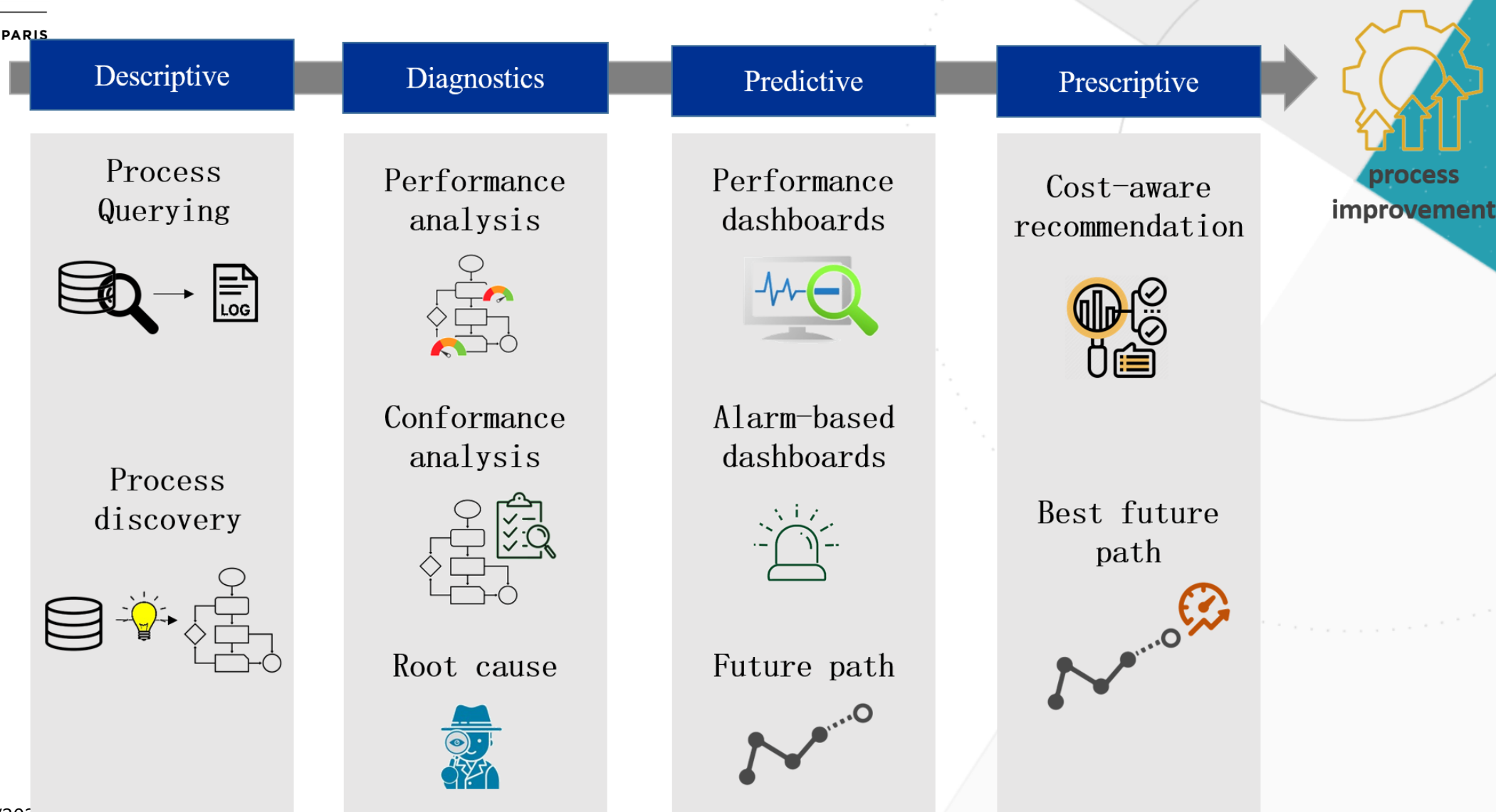
resource

other data

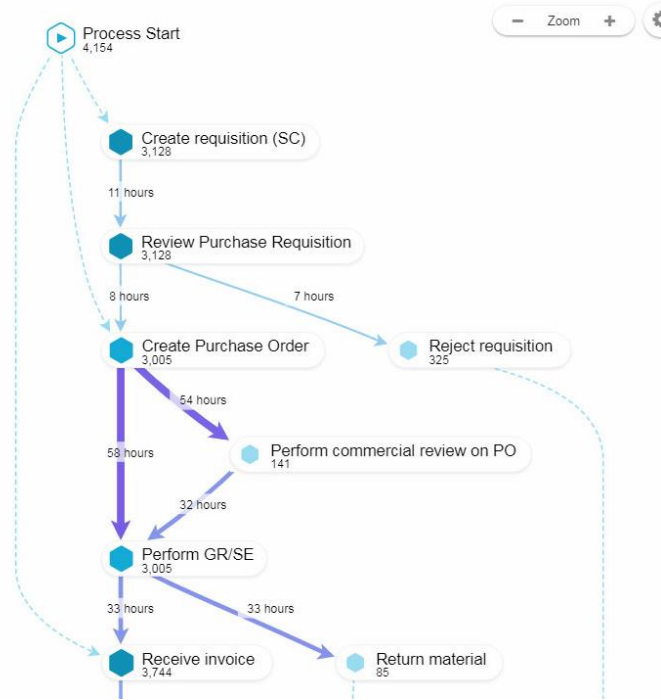
Process Execution data

- **Who (resource) did what (activity) and when (timestamp)**
- **What data, objects, artifacts** are consumed/produced?
- 2 standards: **OCEL**
 - No predefined process instance
 - Events affect objects

Process Mining techniques

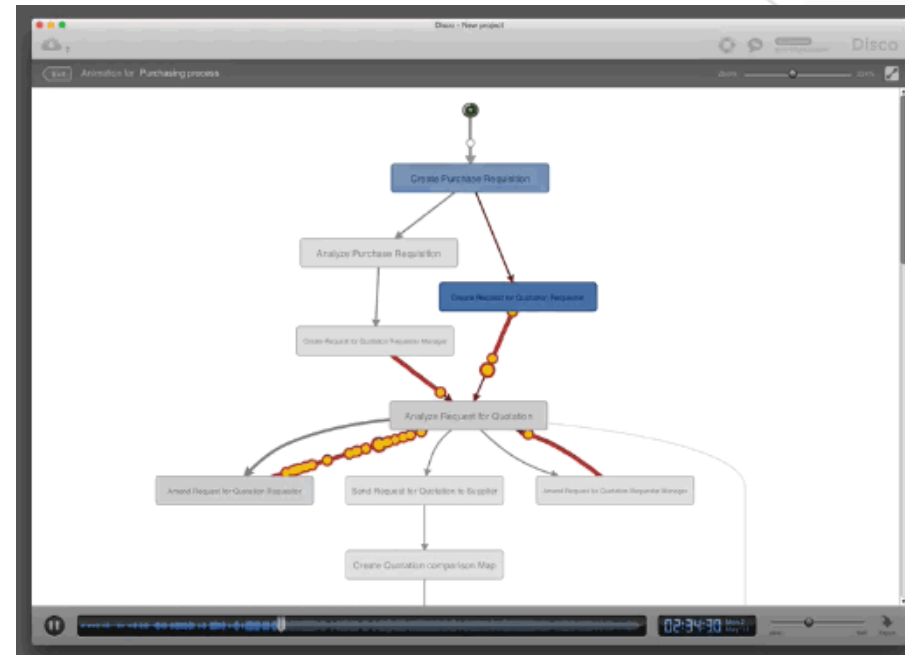
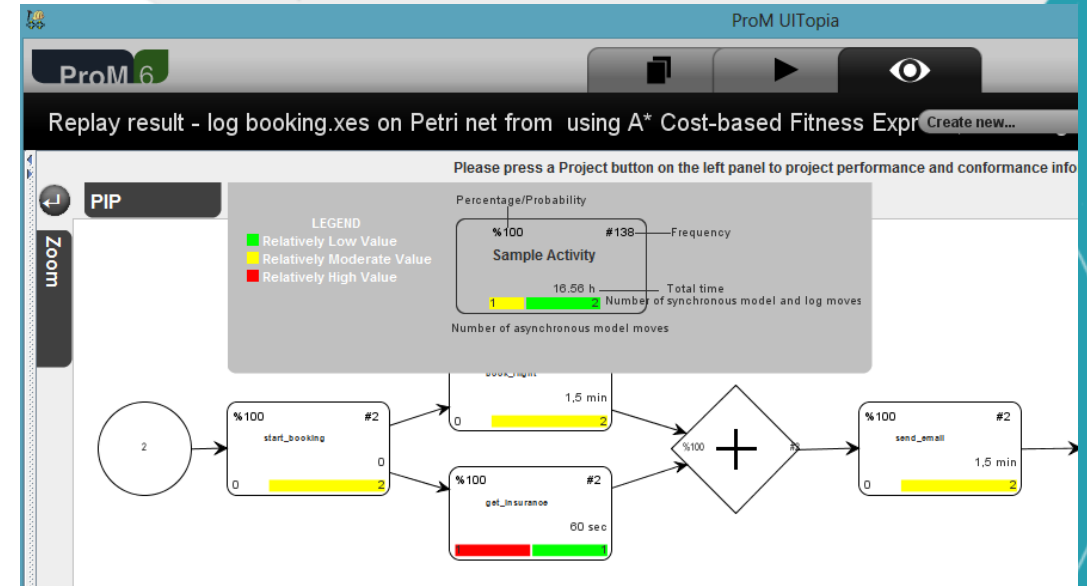


Process discovery & performance analysis



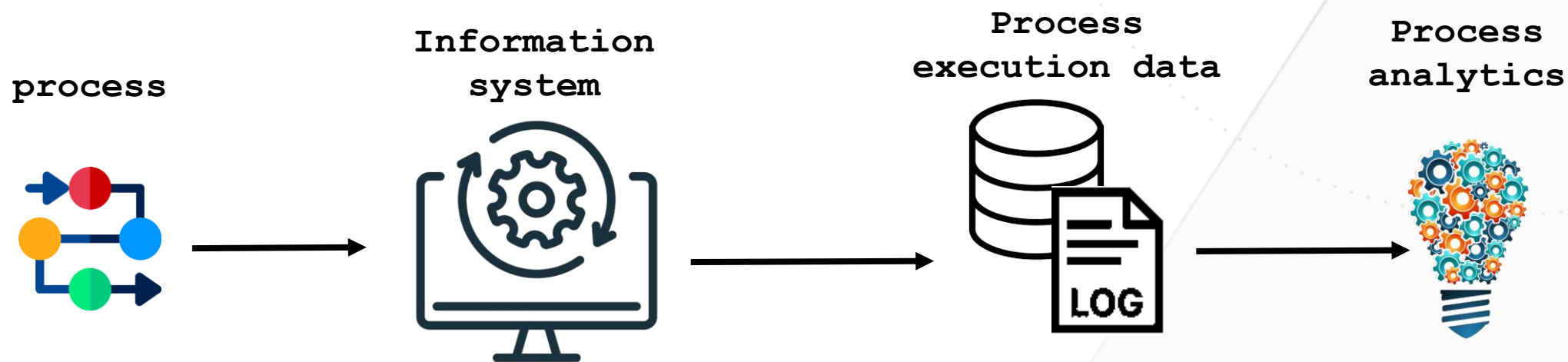
Celonis Raises \$1 Billion
At \$11 Billion Valuation,
Making It New York's —
And Germany's — Most
Valuable Startup

Conformance checking



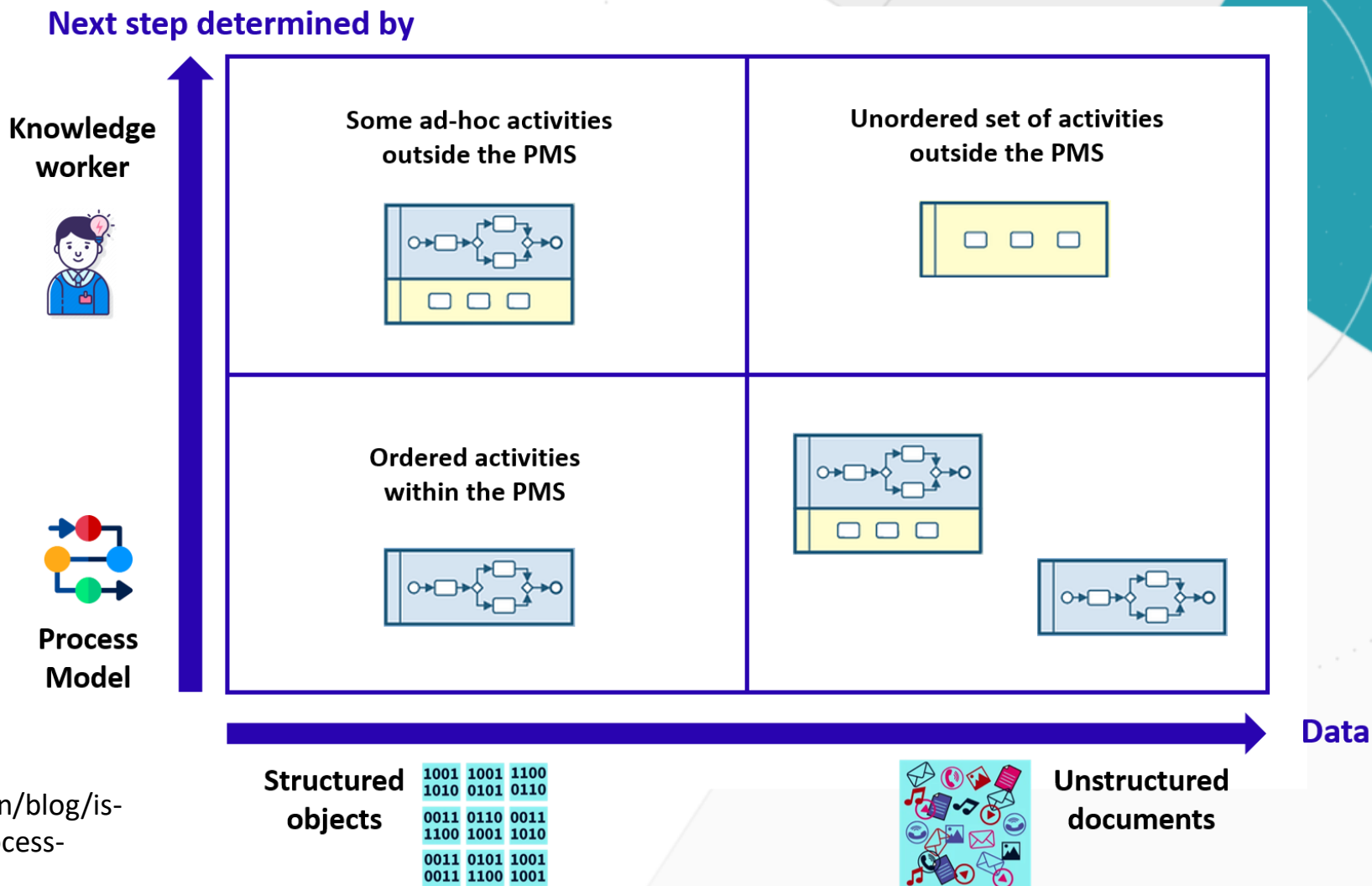
"Idealized" Process Analytics landscape

- Processes are **structured**,
- fully executed within a **process management system**,
- which records **process execution data**,
- that are analyzed using **process analytics** techniques



Spectrum of processes

- from structured to unstructured processes

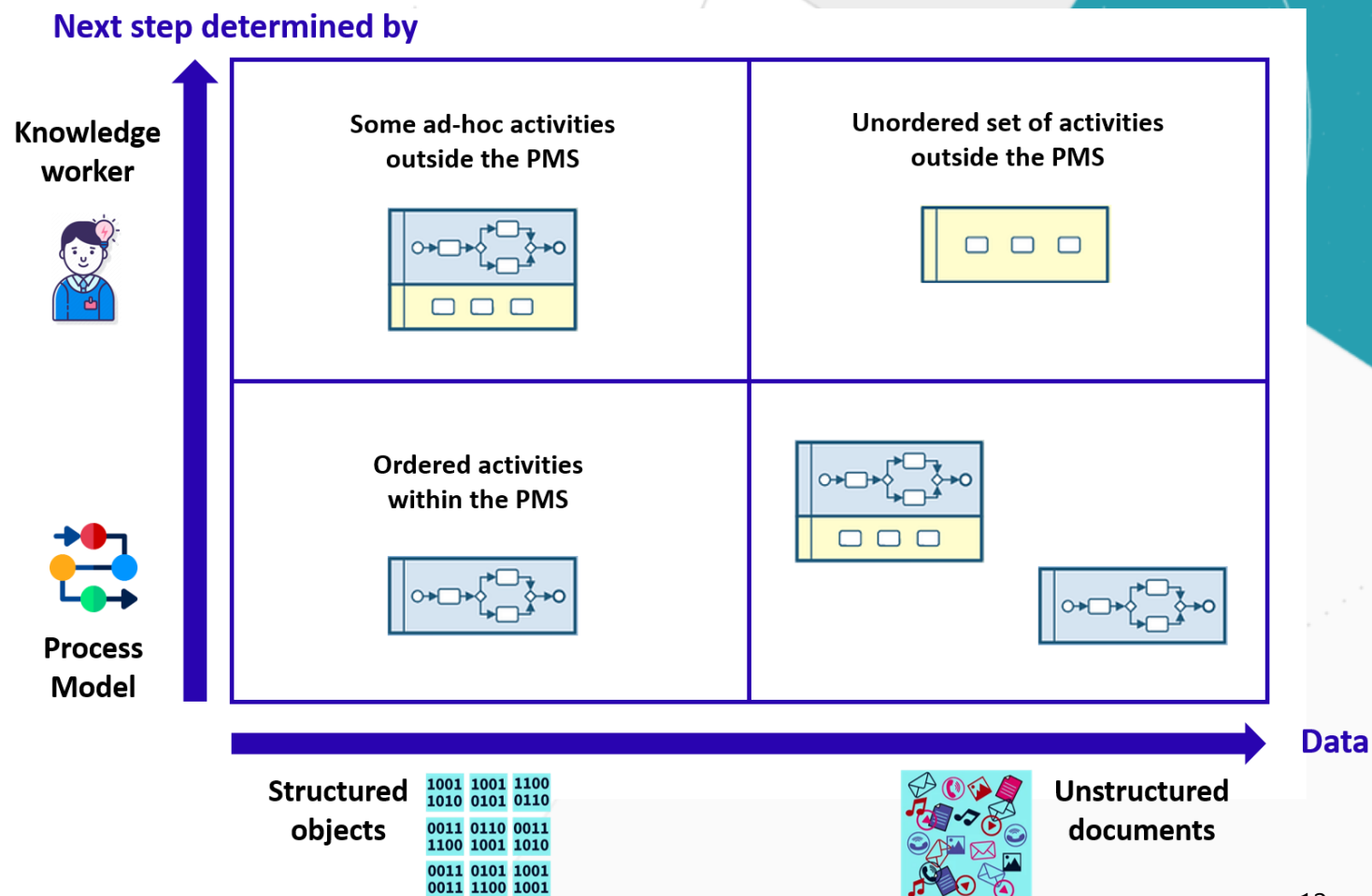


Adapted from: <https://www.nedigital.com/en/blog/is-this-an-ibm-bpm-case-management-or-a-process-management-solution>

Spectrum of processes

workers may:

- **Chat** about processes (e.g. using messaging systems)
- Use **documents/notes** to perform activities
- Use **social medias** (e.g. for recruitment)



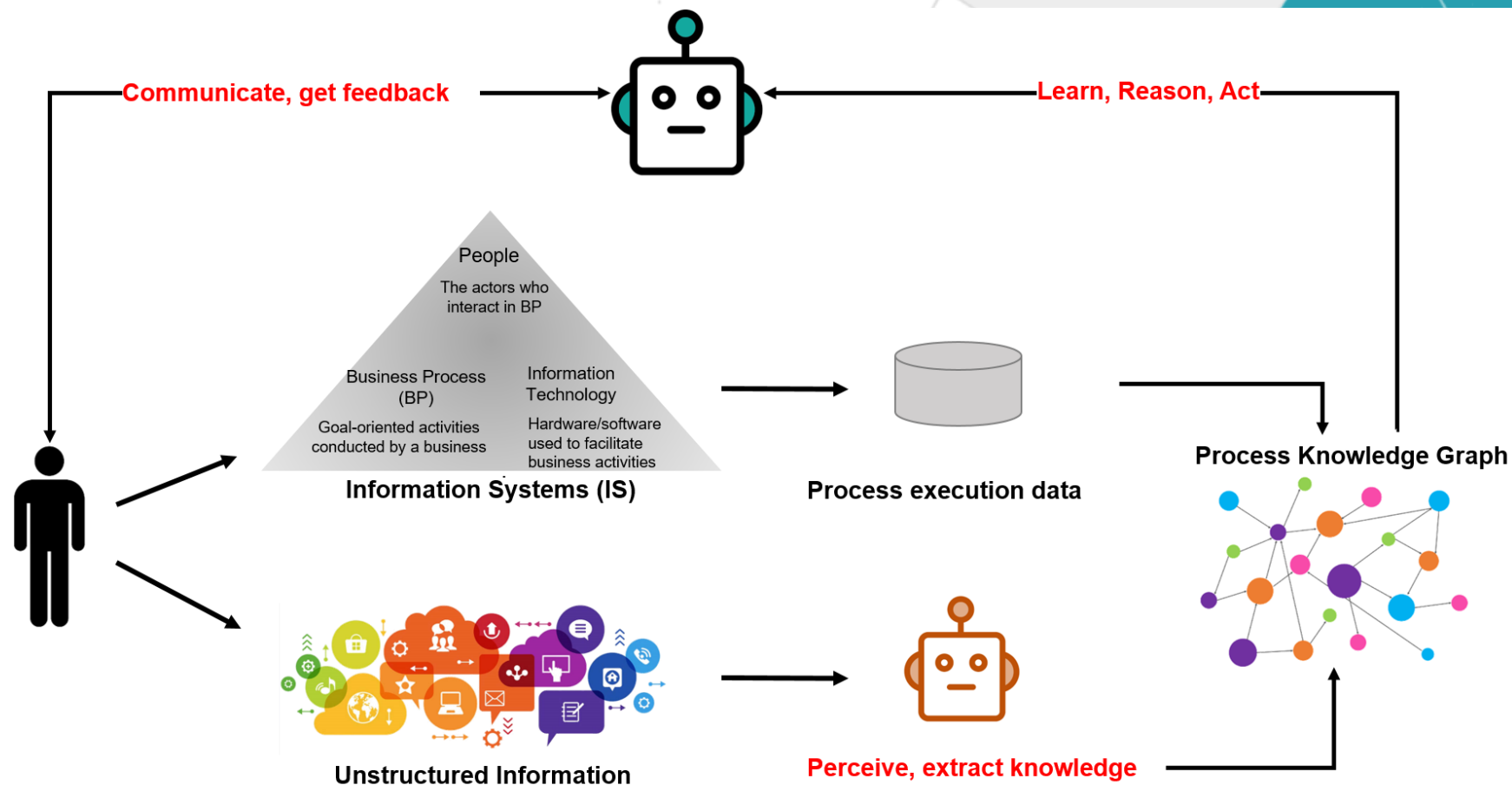
Open Questions

- How to enable **process analytics** on structured and **unstructured data**?
- How to enable process analytics on **unstructured processes**?
- How to extract rich **process information** from **unstructured** data?
- How to **combine** this information with structured process data?
- How to use extracted insights to **assist workers** in their daily tasks?

Towards Cognitive Process Analytics

A system that:

- analyzes structured + unstructured data
- discover, learn, reason
- act to assist workers (reactive + proactive)

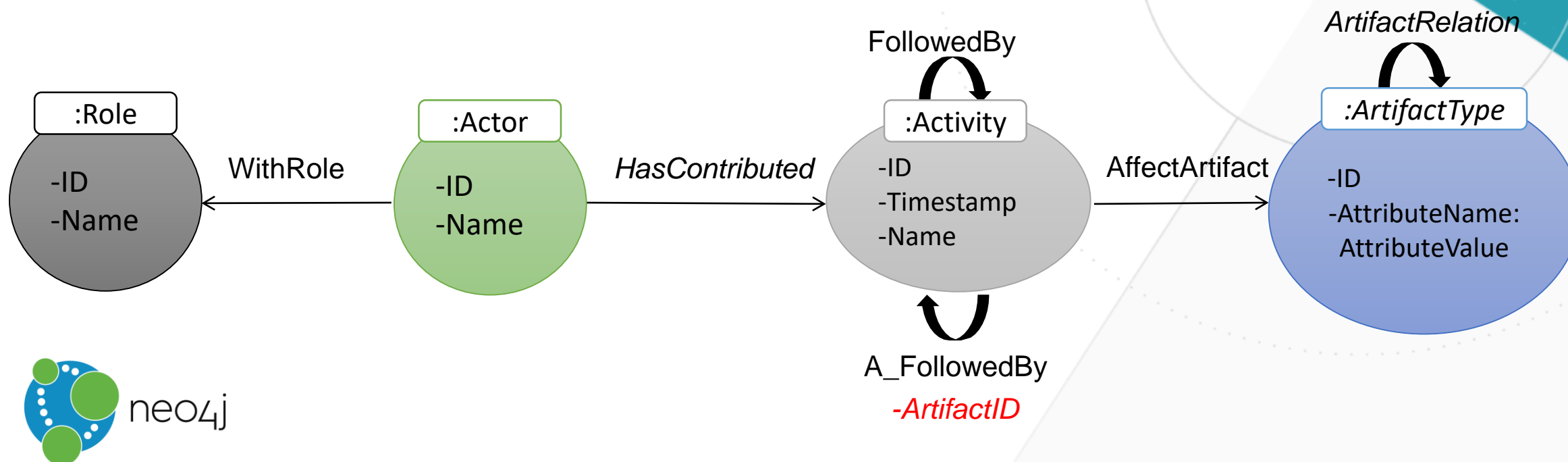


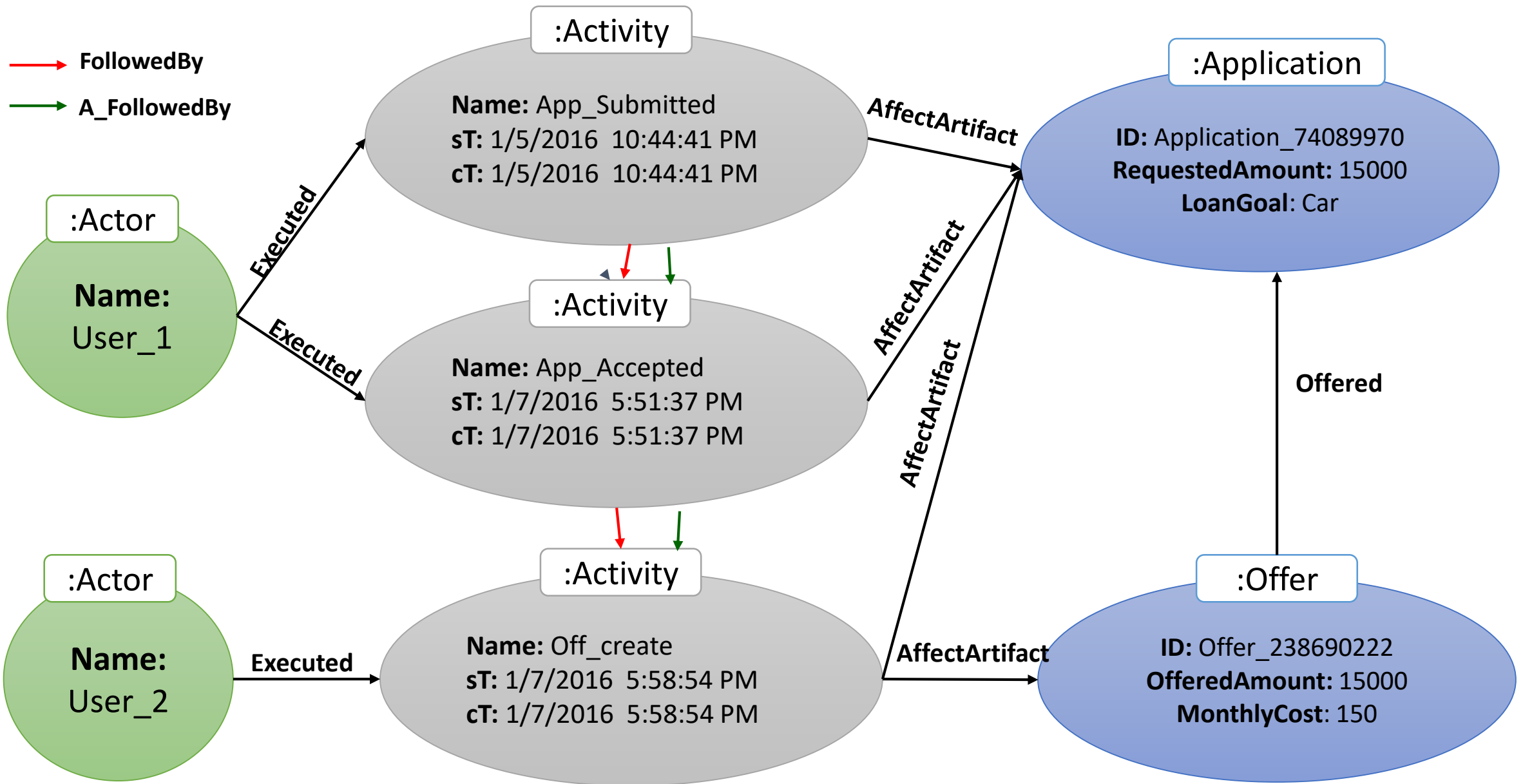
Process Knowledge Graph

- How to effectively store the extracted process information?
- Existing Storage techniques:
 - **Relational databases (SQL like languages):**
 - +Pros:** simple and widely used model
 - Cons:** Querying behavioral aspect is not intuitive (Querying path is complicated with SQL as it requires multiples joins).
 - **Graph databases:**
 - +Pros:** Graph database treats relationship between data as first class citizen
 - Cons:** RDF is not compact (does not support internal structure)

Process Knowledge Graph

- **Graph-based storage using Labeled Property Graph (LPG):** Multi-dimensional event data storage based on OCEL





Conversational Assistant

- How to easily query the stored process data?
- Main limitation of existing process querying approaches:
 - require end users to **learn and master a specific querying language.**
- Process querying is primarily **business-driven** and should be accessible to business experts who may be **inexperienced in database querying.**

Agenda

- Introduction to Processes & Process Analytics
- Spectrum of processes
 - From structured to unstructured processes
- Conceptual models to enable Cognitive Process Analytics
 - Process Knowledge Graph
 - Conversational Assistants
- Domain specific models: Emailing Systems use case
 - Discover Process Knowledge Graph from emails
 - Conversational Assistant to query Process Knowledge Graph

Research Problem

- How to enable cognitive process analysis on email data?

1. How to extract process related information from emails?

activities, actors, business artifacts, etc.



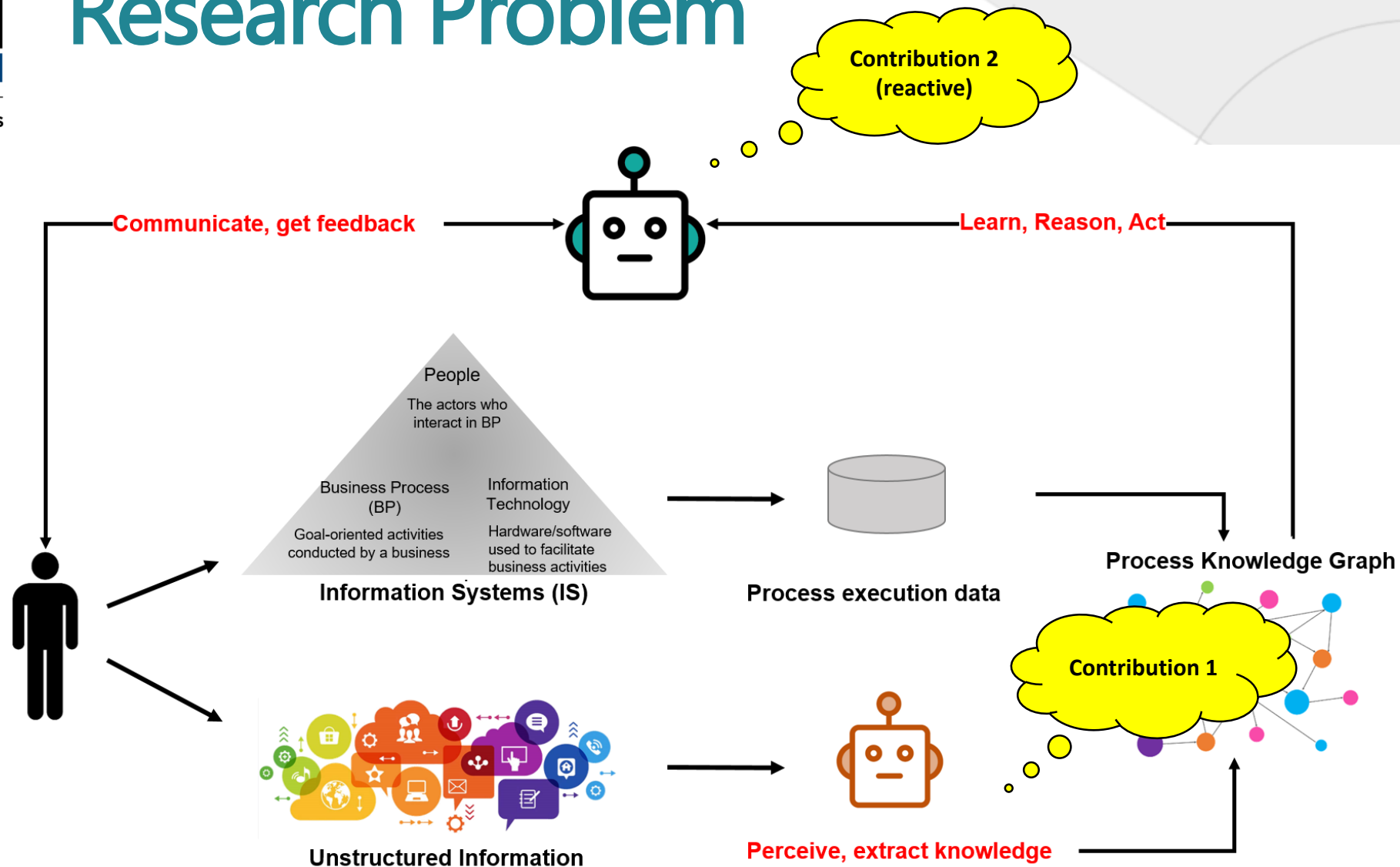
3. How to assist end users in querying the process data?

Querying the process data should be accessible to business experts

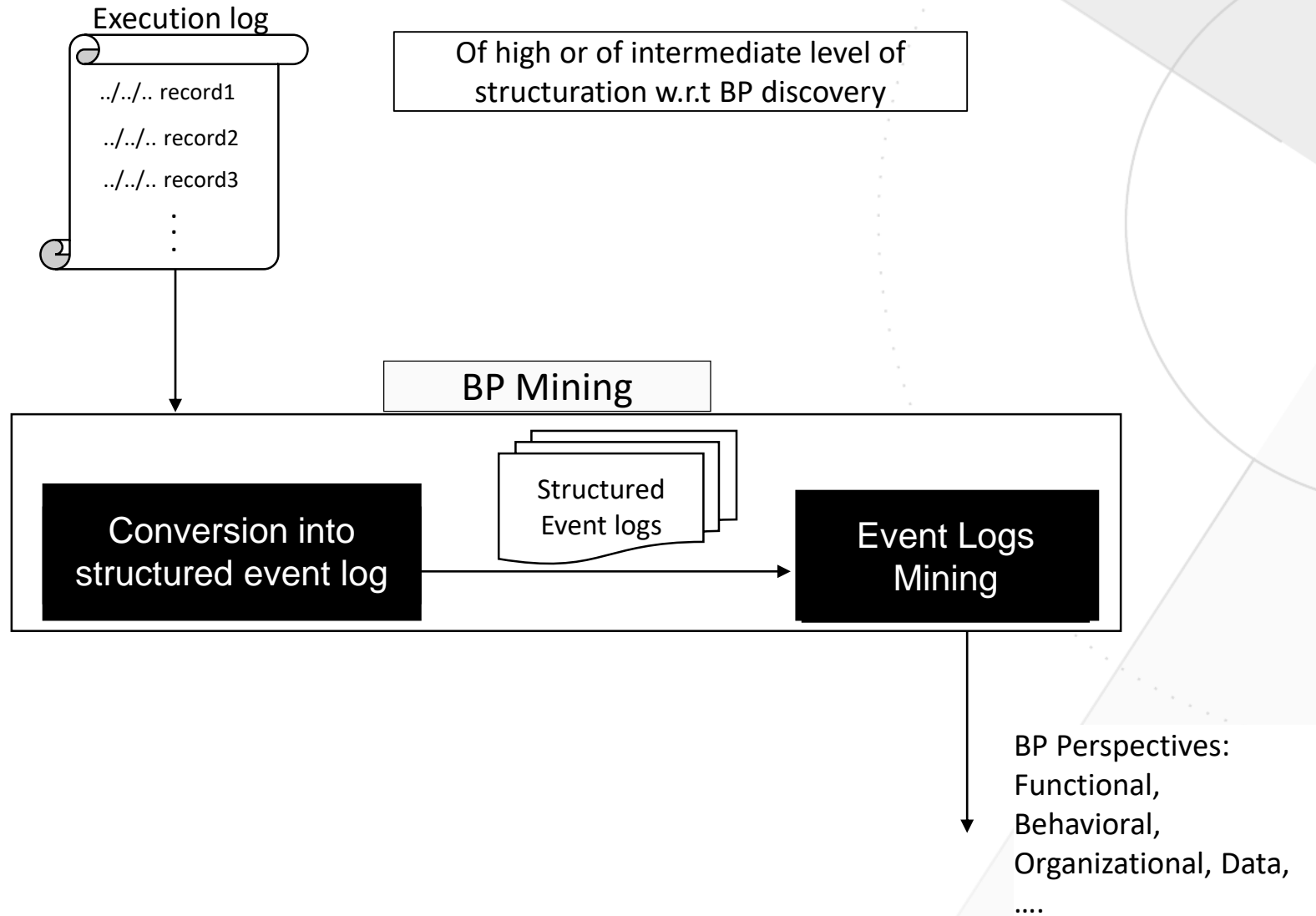
Contribution 1

Contribution 2

Research Problem



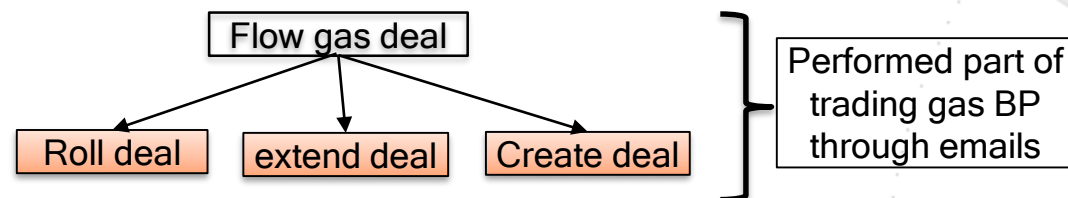
PROCESS MINING



Emailing systems for performing BP activities



EMAILING SYSTEMS FOR PERFORMING BP ACTIVITIES



Date: Wed, 15 Nov 2000 02:24:00 -0800
From: daren.farmer@enron.com
To: aimee.lannou@enron.com
cc:
Subject: Re: Flow w/ no nom

Rolled deal 454057 to cover flow at mtr 5192.

d

Aimee Lannou 11/10/2000 02:17 PM
To: Daren J Farmer/HOU/ECT@ECT
cc:
Subject: Flow w/ no nom

Meter 1601 last deal 412219 for 10/00 flowed 11/9

Meter 5192 last deal 454057 for 10/00. flowed 11/3-4

Email1 From Enron database

Informing about rolling a deal

Informing about extending a deal

conversation history

Informing about gas flowing of two deals

Date: Tue, 9 Jan 2001 06:48:00 -0800 (PST)
From: daren.farmer@enron.com
To: aimee.lannou@enron.com
cc:
Subject: Re: Dec 00

I extended 454057 for the month of December.

D
Aimee Lannou 01/09/2001 12:43 PM

To: Daren J Farmer/HOU/ECT@ECT

cc: Edward Terry/HOU/ECT@ECT

Subject: Dec 00

Daren - meter 5192 flowed 8 dth on 12/19, 33 dth on 12/20 and 2 dth on 12/29. The last deal for this meter was 454057 in Nov 00. Could you please extend this deal for these 3 days or create a new one? Please let me know.

AL

Email2 From Enron database

➔ The importance of analyzing emails for covering the traces of:

- BP that are totally or partially performed outside structured information systems
- Employees interactions regarding BP activities

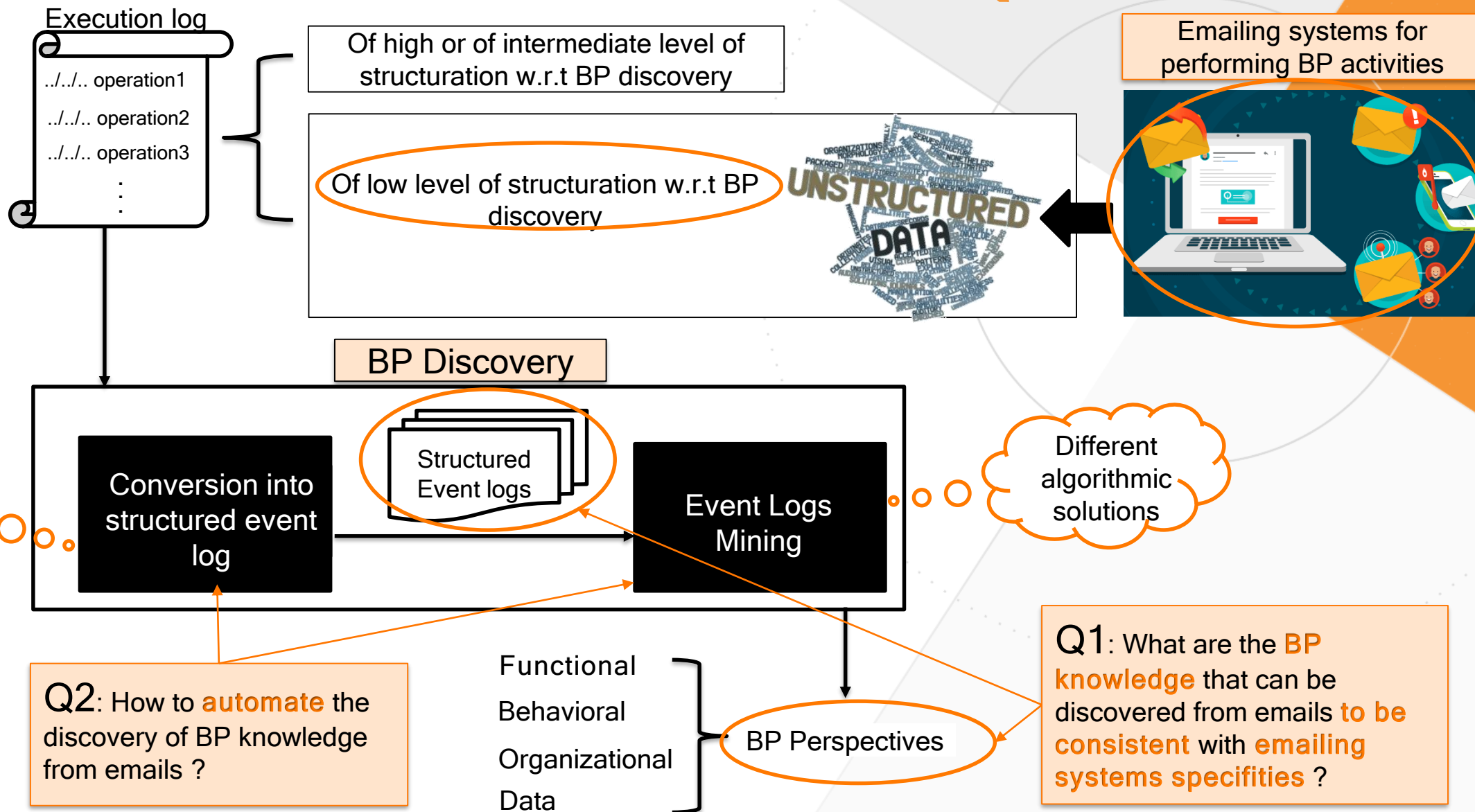
➔ Enhance Business Process Management (BPM) by covering unmanaged and less structured BP traces

Informing about gas flowing

Requesting extending deal

Requesting creating new deal

GENERAL CONTEXT & MAIN RESEARCH QUESTIONS



Main Objectives

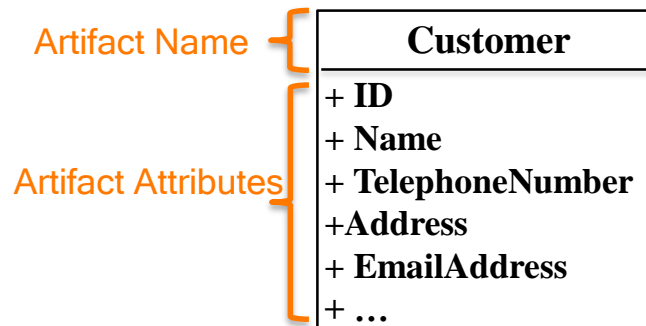
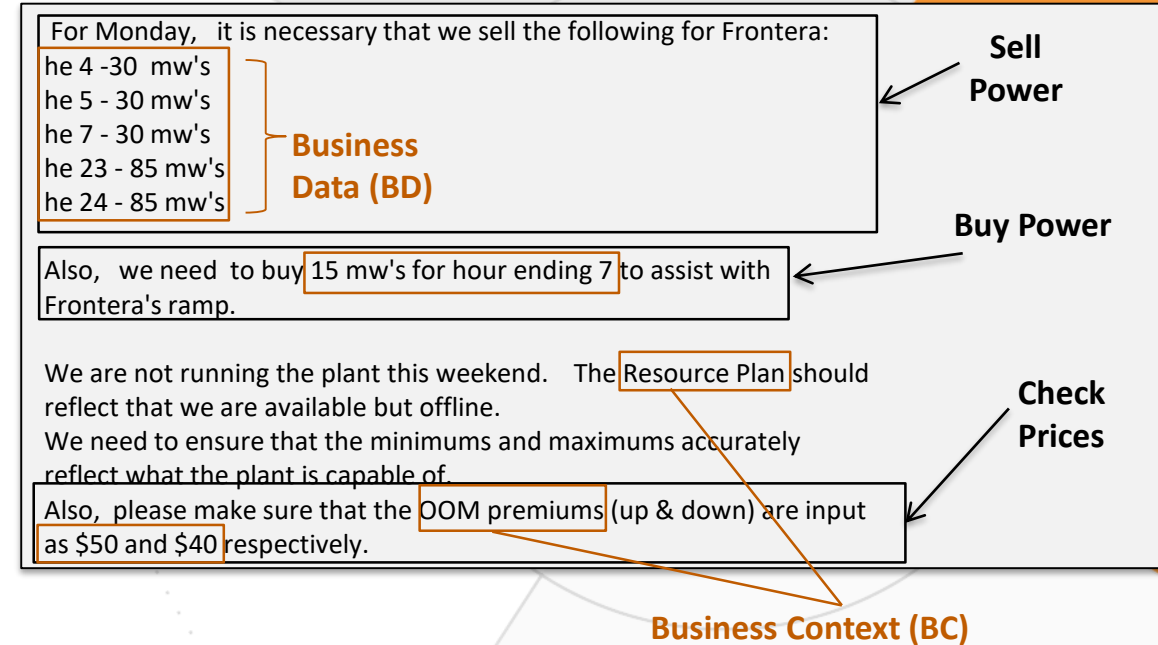
Objective 1: Formalize the definition of BP knowledge (i.e. Multiple BP perspectives & event log structure) to be discovered from emails

Objective 2: Propose a totally unsupervised approach for discovering BP w.r.t multiple perspectives

- Without requiring a priori BP information & considerable human intervention
- While allowing the discovery of multiple activities per one email

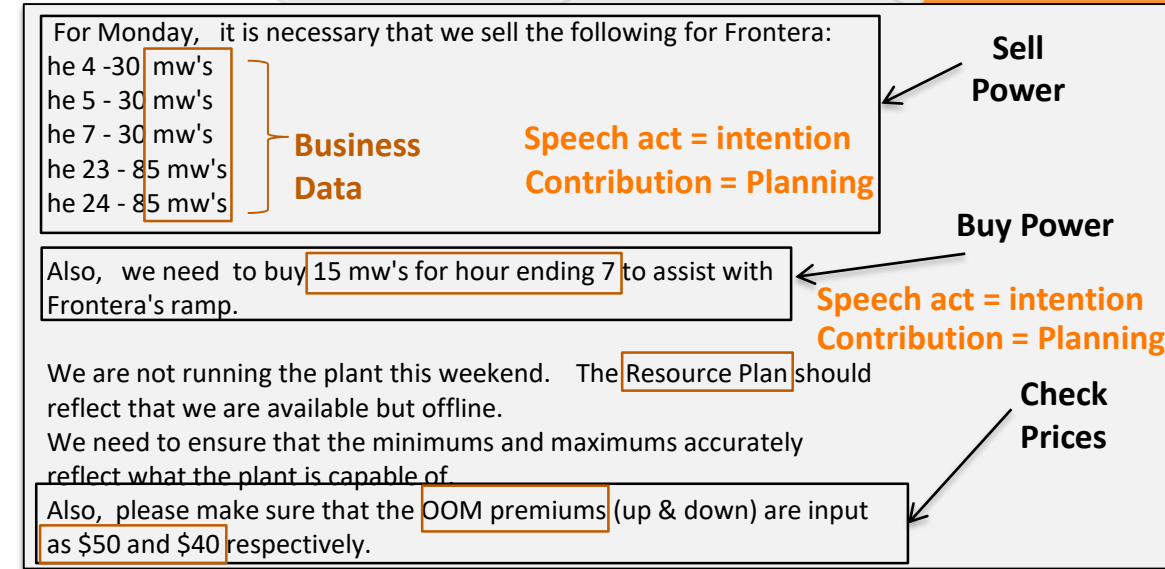
OBJECTIVE 1: BP KNOWLEDGE & EMAILS

- **BP Activity:** Activity Name + Business Information
- **Functional perspective:** Adopt the notion of **BP fragment** rather than complete BP due the incompleteness of BP traces in emails
- **Data Perspective:** Adopt **artifact** centric approach in the context of emails

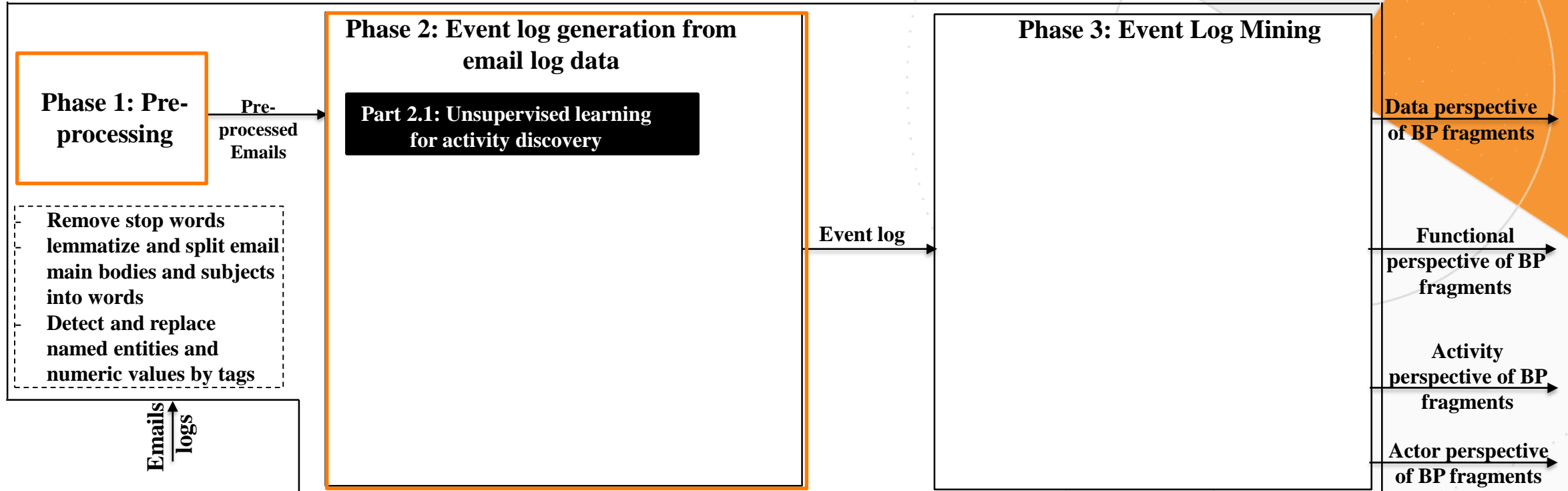


OBJECTIVE 1: BP KNOWLEDGE & EMAILS

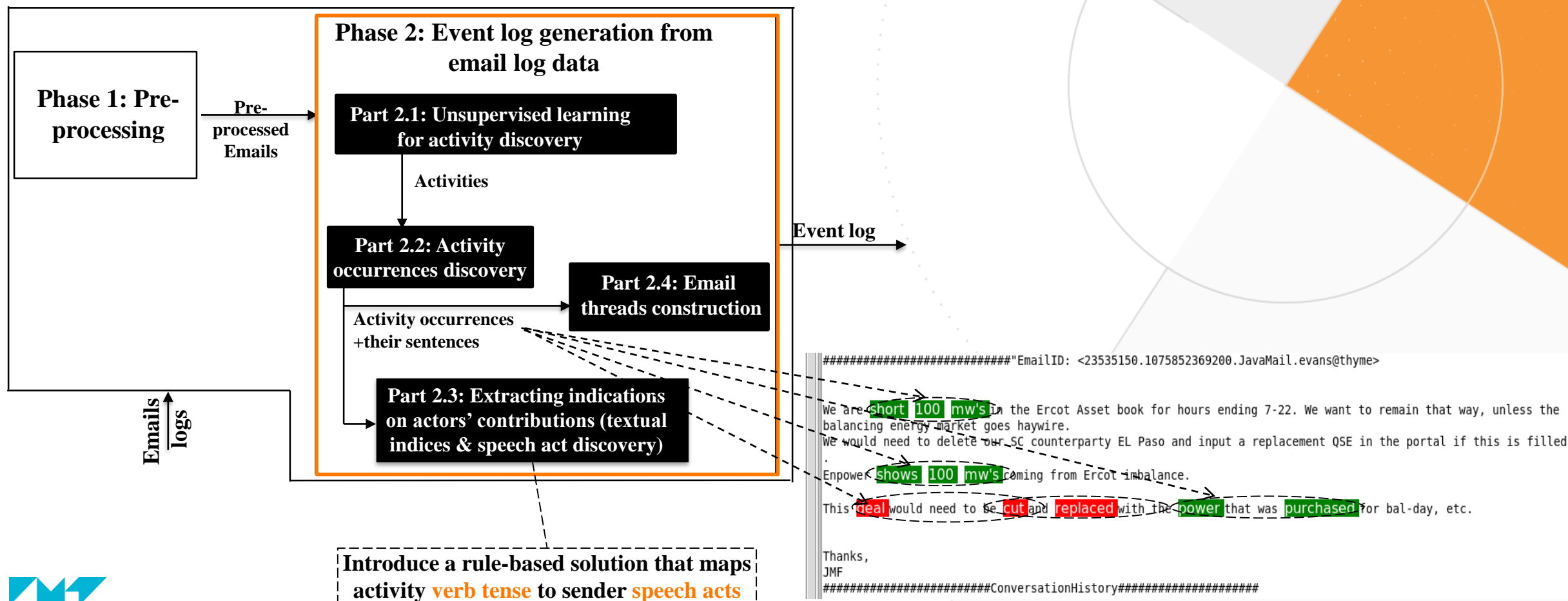
- **BP Activity:** Activity Name + Business Information
- **Functional perspective:** Adopt the notion of **BP fragment** rather than complete BP due the incompleteness of BP traces in emails
- **Data Perspective:** Adopt **artifact** centric approach in the context of emails
- **Behavioral Perspective:**
 - Adopt declarative approach for activity control flow;
 - Introduce new event type that combines activity & speech act;
 - Information act
 - Intention act
 - Request act
 - Request information act
- **Organizational perspective:** Consider **multiple actors with various contributions (not limited to executing activities)** when performing an activity
 - Execution
 - Information
 - Planning
 - Request
 - Request information
 - Observation



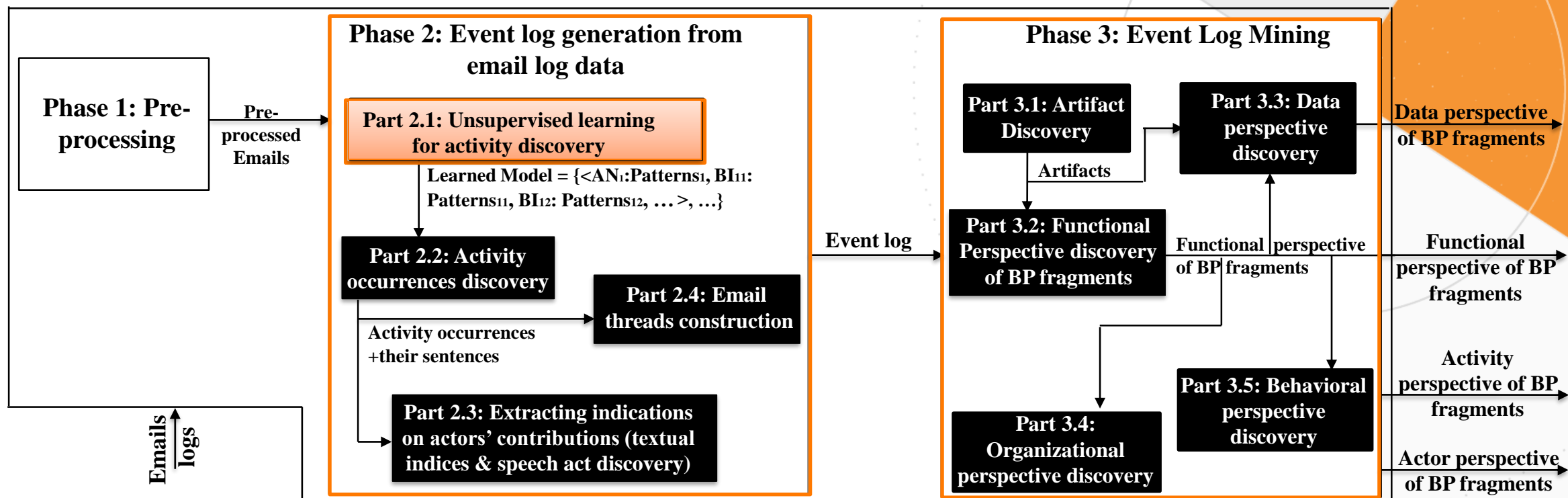
OBJECTIVE 2: A TOTALLY UNSUPERVISED APPROACH FOR BP KNOWLEDGE DISCOVERY



OBJECTIVE 2: A TOTALLY UNSUPERVISED APPROACH FOR BP KNOWLEDGE DISCOVERY



OBJECTIVE 2 : A TOTALLY UNSUPERVISED APPROACH FOR BP KNOWLEDGE DISCOVERY



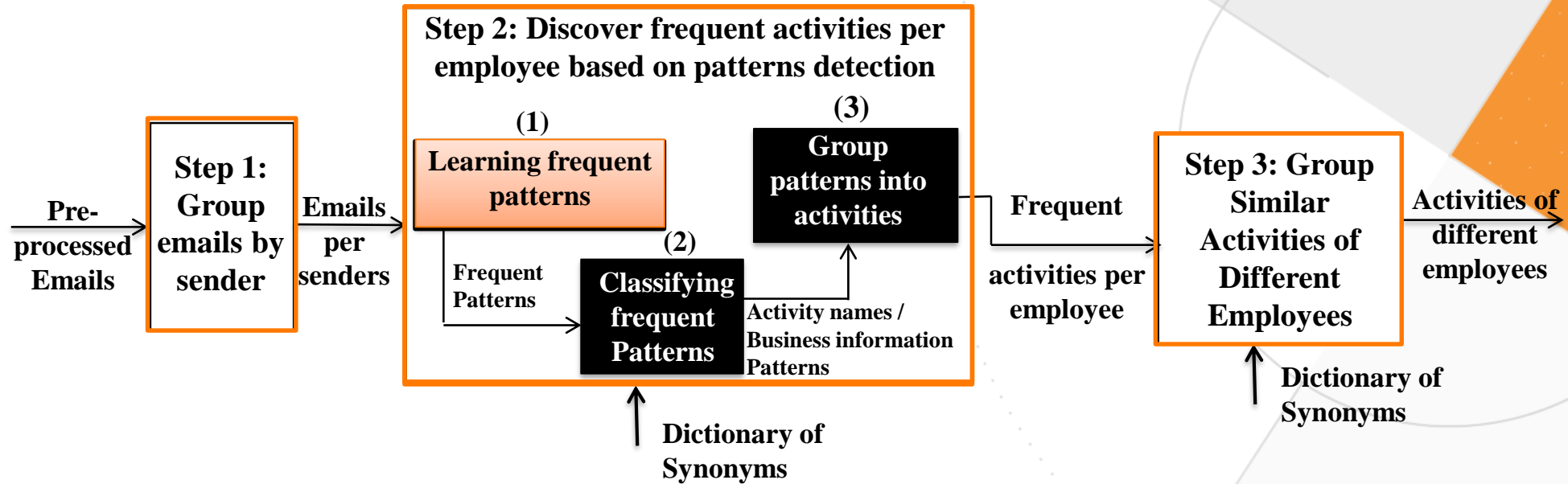
Main Assumptions:

- 1 For each BP, \exists a set of actors that frequently execute the same activities
- 2 One employee that frequently expresses an activity in emails would use close expressions

→ Frequent patterns of words



UNSUPERVISED LEARNING FOR ACTIVITY DISCOVERY: MAIN STEPS



DISCOVERING PATTERNS: MAIN PROPOSITIONS

- 1) Consider **low dispersion constraints** when discovering frequent patterns of words to avoid non-significant patterns

Example of low dispersed pattern in a text

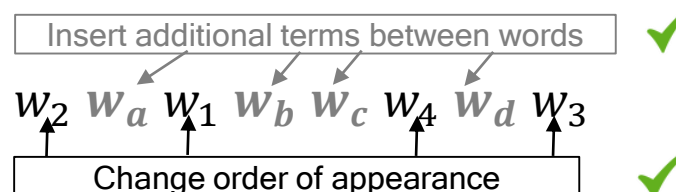
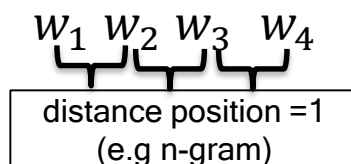
I have **created** deal **tickets** 241558 and 241560 for July 99 - March 00 based on the info below. Due to time constraints, I have not researched pricing and volumes, I trust that the info supplied to me agrees with the contract.

Example of highly dispersed pattern in a text

I have **created** a spreadsheet to assist in the tracking and booking of the gas supply related to Tenaska IV. If you have suggestions or comments on the file, contact either me or Mark.

Megan - You can copy the Jan 01 Est tab and update the volumes with actuals. The file should calculate the resulting settlement with Tenaska. I will then update the demand fee on the deal **ticket** to true-up the month. Since we haven't gone through actuals in the spreadsheet, we may have to do some tweaking with the formulas. Anyway, this should give us a good start and you will be able to see where we closed the month in Logisitics.

- 2) Do not impose constraints concerning the **order of appearance of words** (\neq Sequential Pattern Mining algorithms, e.g. n-gram)



- 3) Introduce the notion of **patterns of concepts** to tolerate the use of synonyms rather than identical words in emails

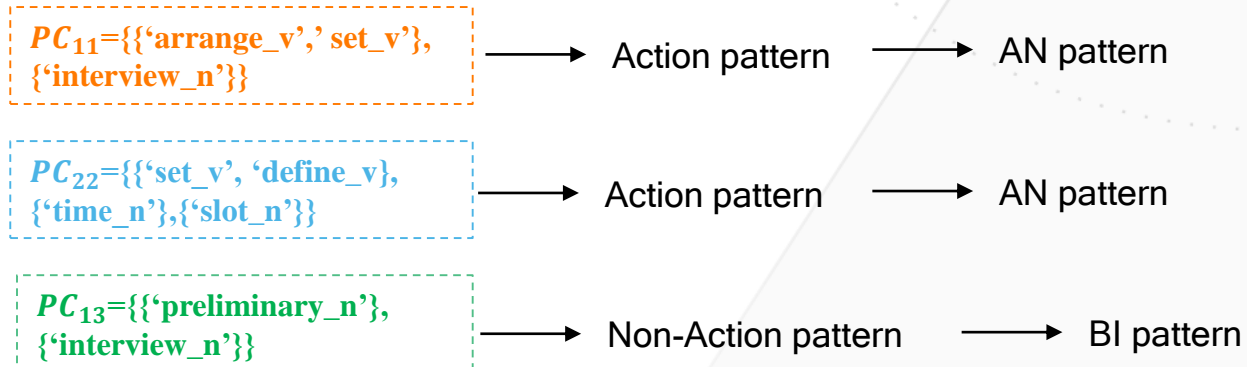
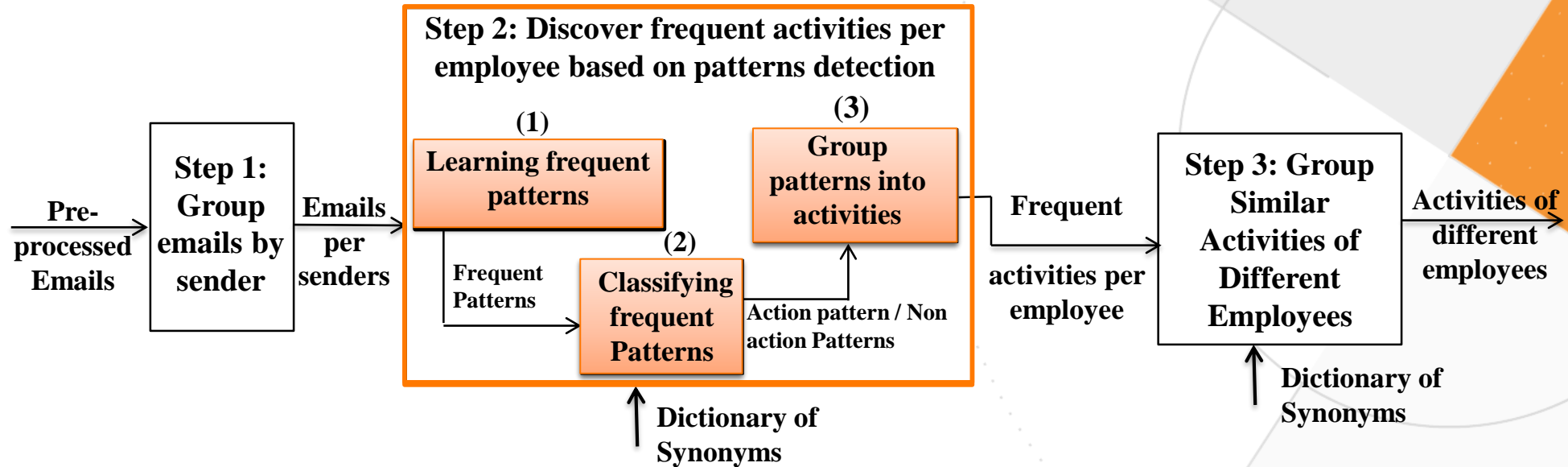
Pattern of concepts = $\{C1, C2, C3, \dots\}$ where $C_i = \{\text{syn1}_i, \text{syn2}_i, \dots\}$ such that all words **have pairwise reciprocal synonymy relations to potentially refer to a unique meaning**

email1 = ['purchase', 'power', ...]

email2 = ['buy', 'power', ...]

\Rightarrow In common Pattern of concept
 $\{ \{ \text{'purchase'}, \text{'buy'} \}, \{ \text{'power'} \} \}$
 C1 C2

UNSUPERVISED LEARNING FOR ACTIVITY DISCOVERY: APPROACH OVERVIEW

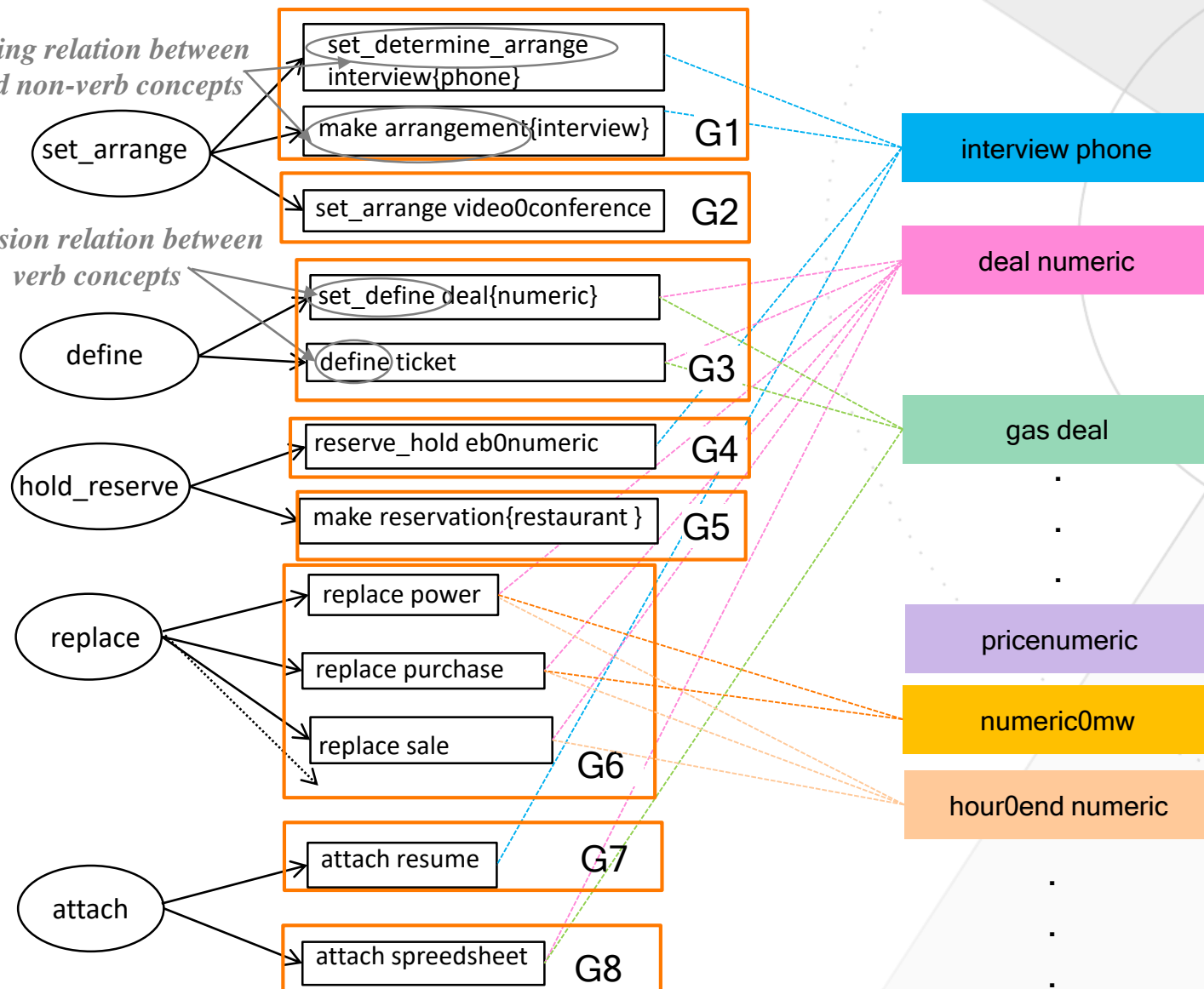


(1) Organize Action patterns by similar realization type

(2) Identify Highly correlated BI patterns to the action patterns (In term of coexistence in the same email)

Rewording relation between verb and non-verb concepts

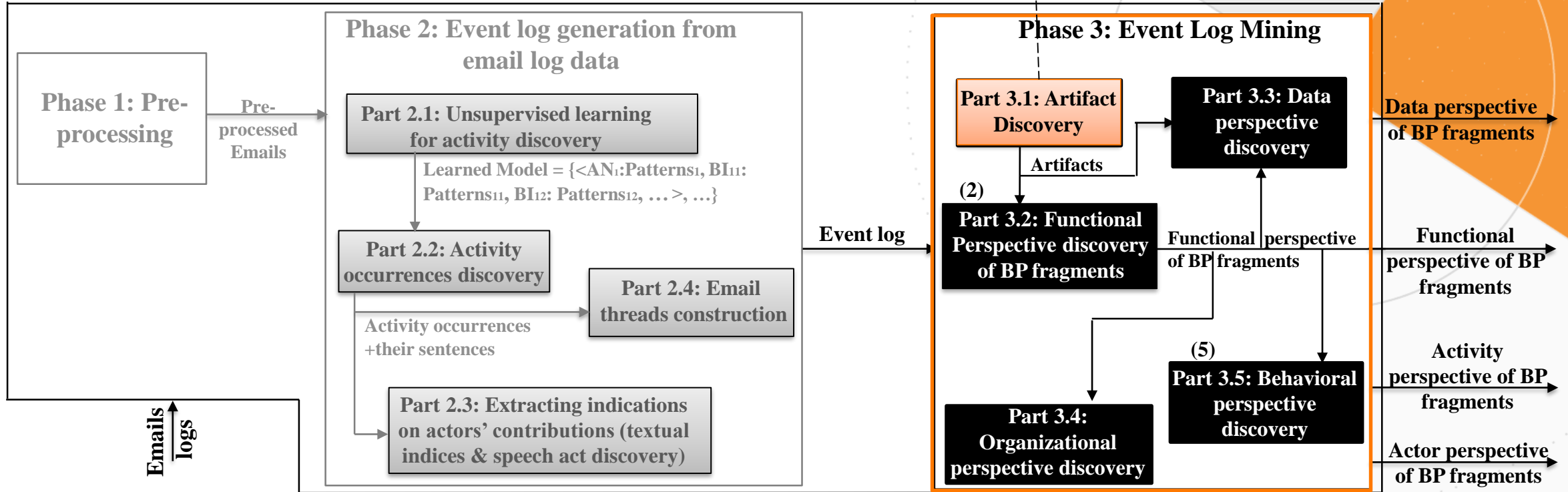
Inclusion relation between verb concepts



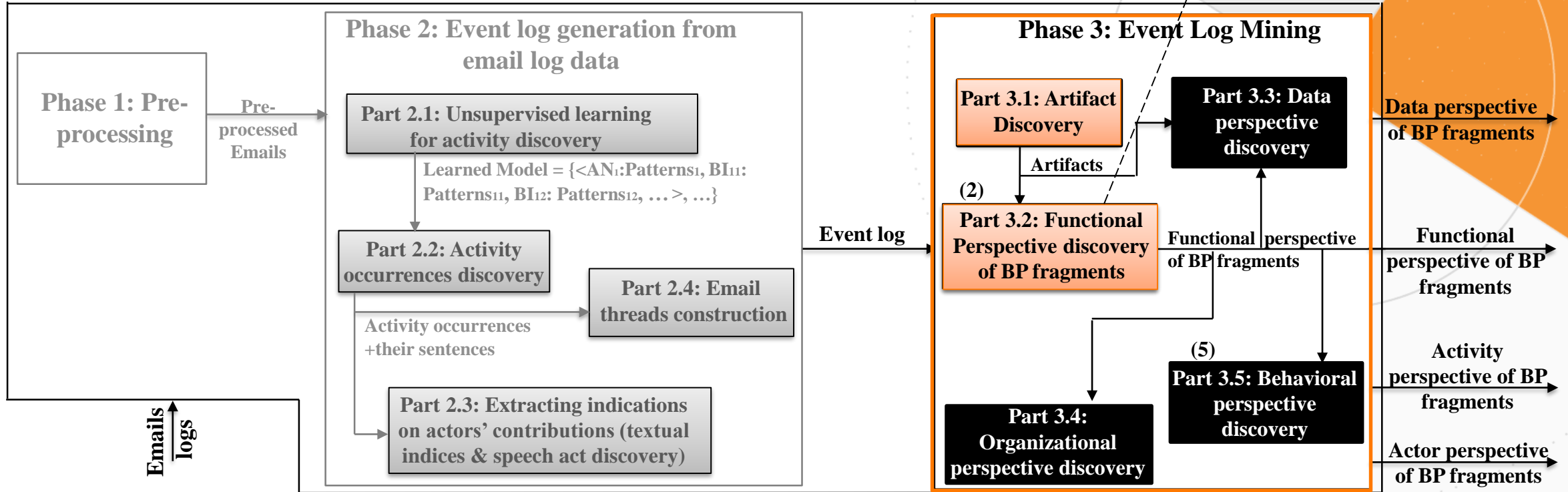
(3) Obtain activities:
Regroup action patterns
of the same realization
type by business context
similarity

EVENT LOG MINING: APPROACH OVERVIEW

One activity could manipulate multiple artifacts



EVENT LOG MINING: APPROACH OVERVIEW



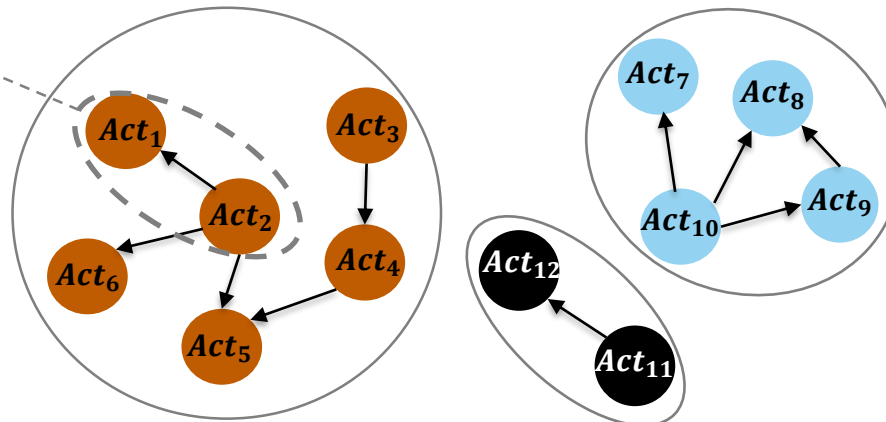
BP FRAGMENT DISCOVERY: MAIN ASSUMPTIONS & PROPOSITIONS

Assumption 1 : Each BP fragment differs from other ones by a set of causality relations between BP elements (i.e. activities, actors & artifacts)

Main Proposition: Discover causality relations between BP elements as a first step towards grouping activities into BP fragments

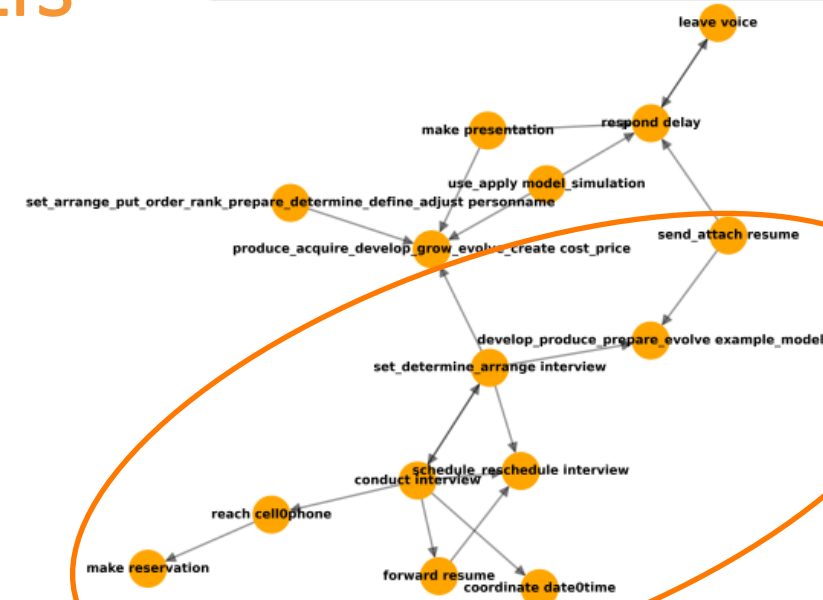
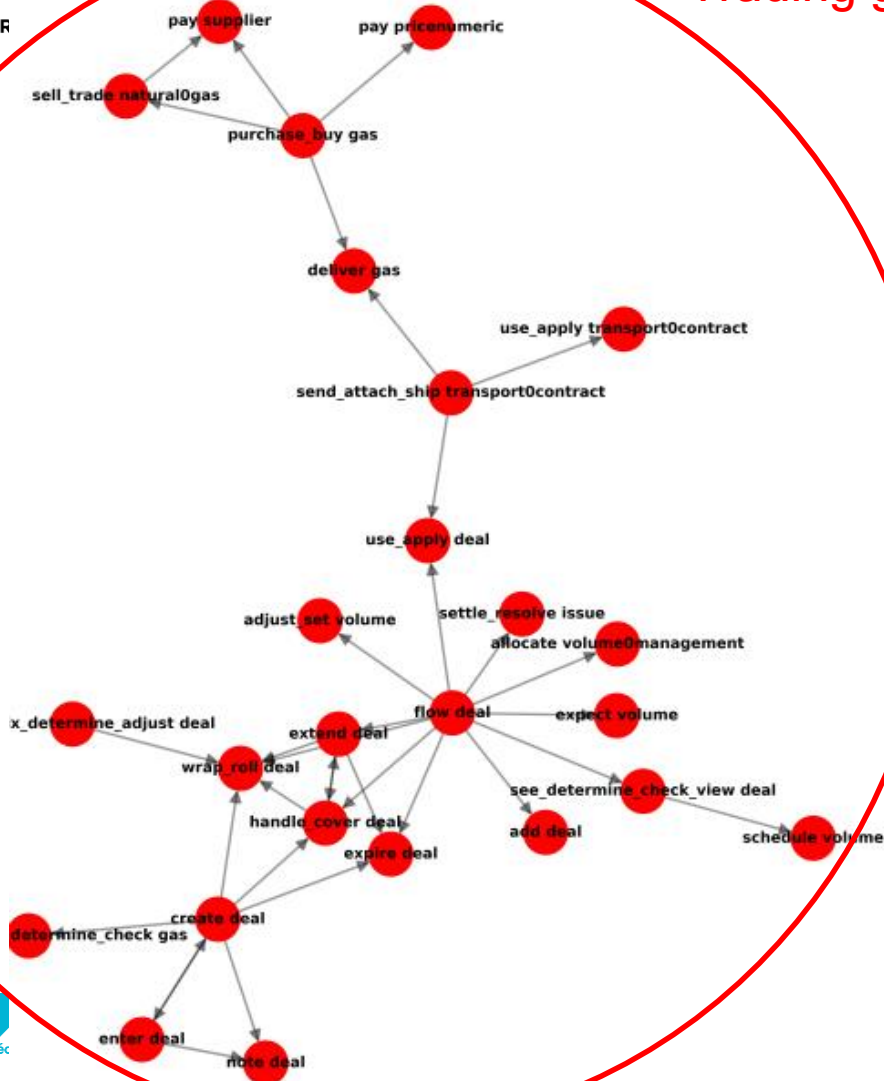


Act_i & Act_j must have causality relations referring to common BC and actor groups



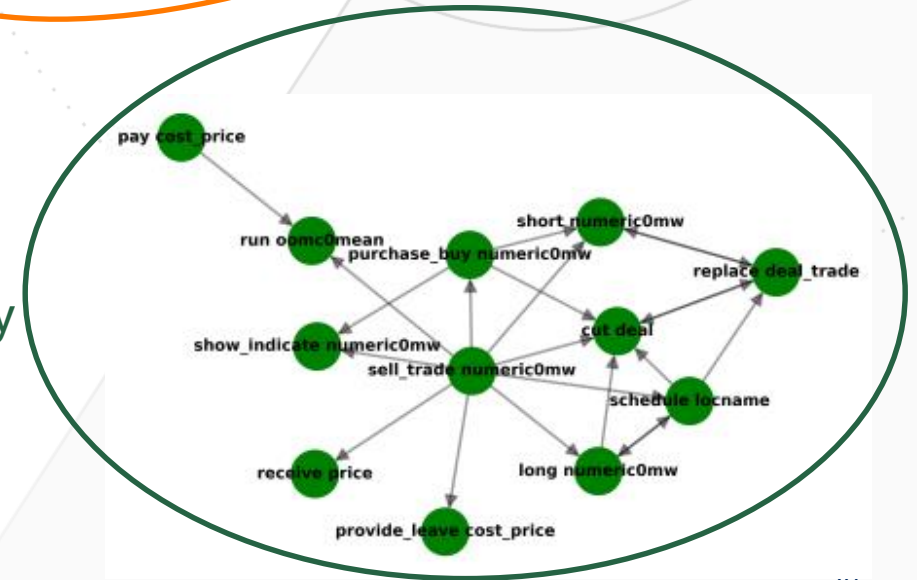
BP FRAGMENT DISCOVERY: RESULTS

Trading gas

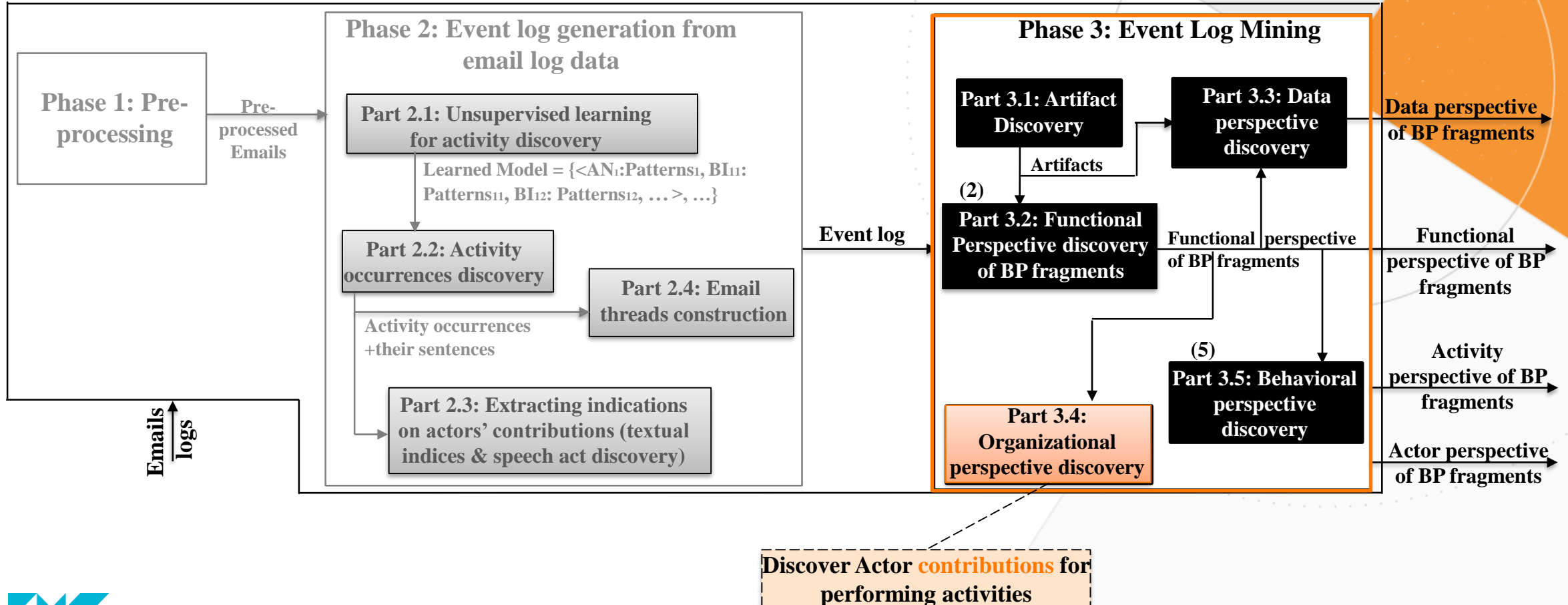


Organizing Interviews

Trading Electricity Power

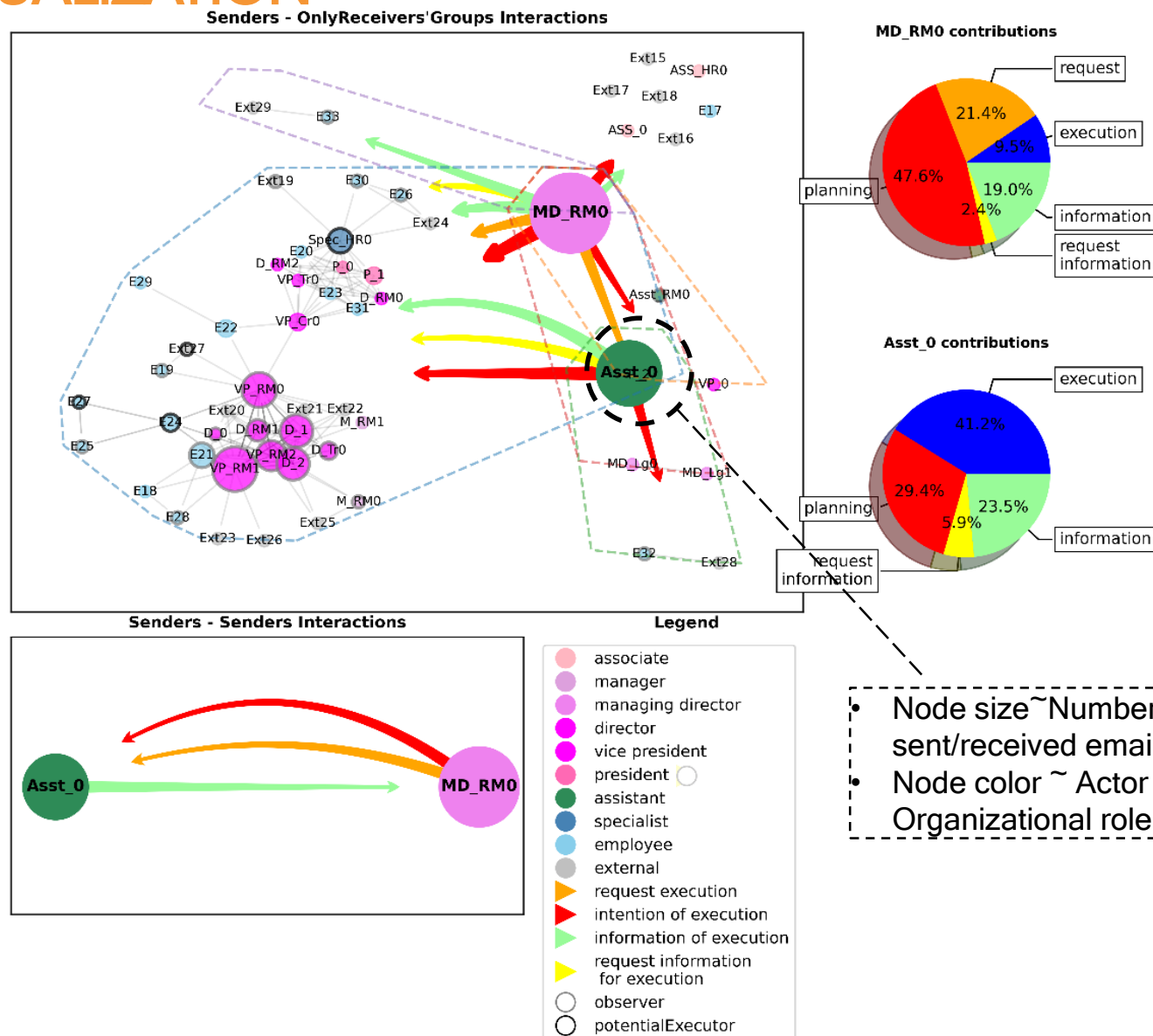


EVENT LOG MINING: APPROACH OVERVIEW

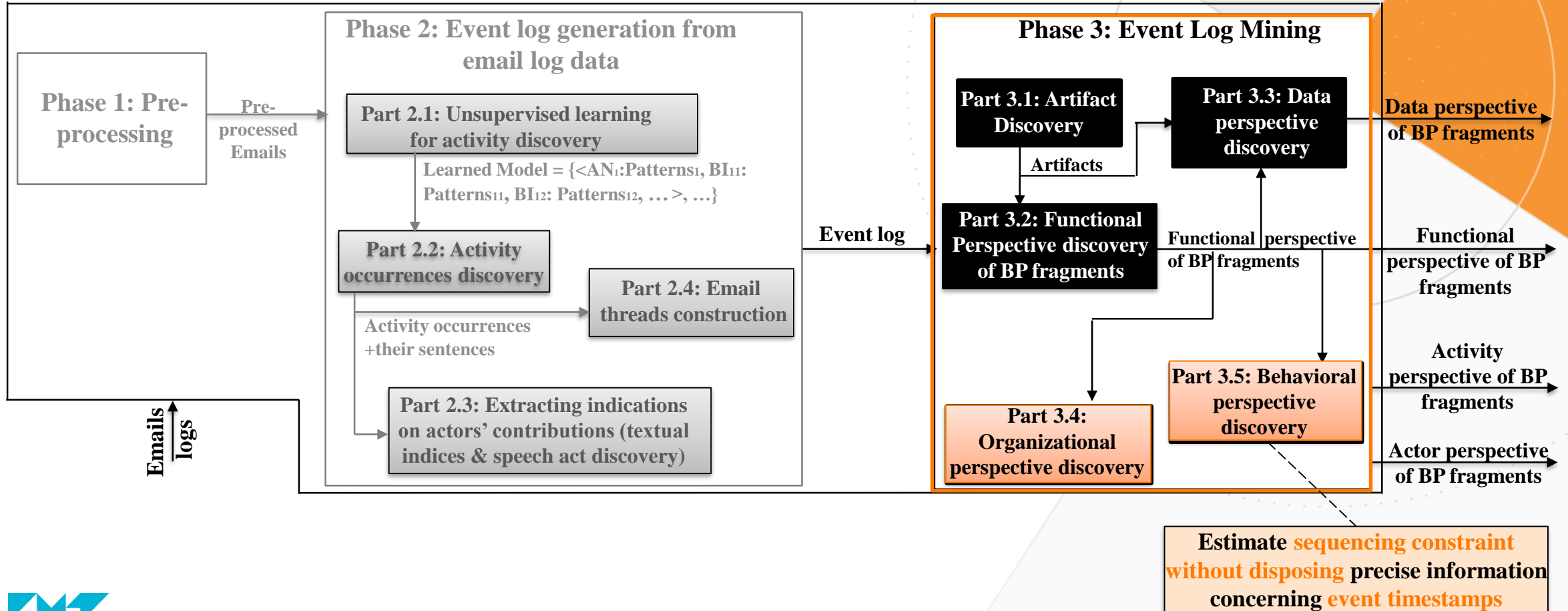


ORGANIZATIONAL PERSPECTIVE DISCOVERY: EXAMPLE OF VISUALIZATION

Organizational perspective of the
activity set_determine_arrange
interview

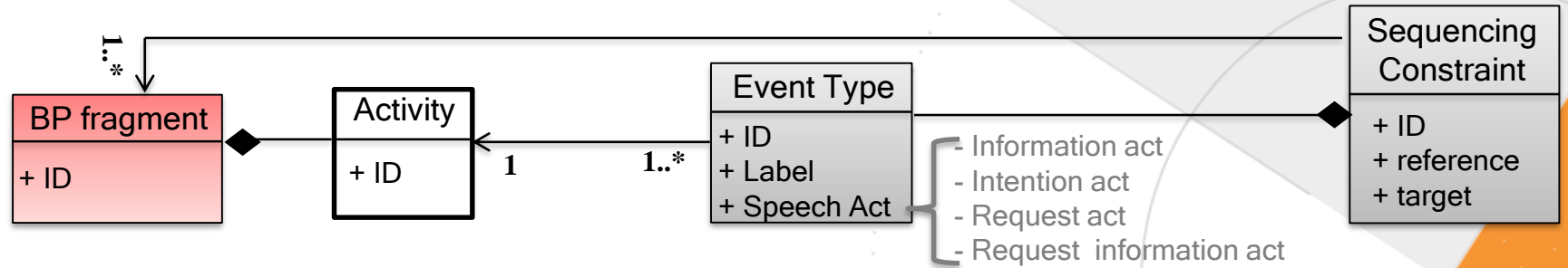


EVENT LOG MINING: APPROACH OVERVIEW



BEHAVIORAL PERSPECTIVE: DEFINITION

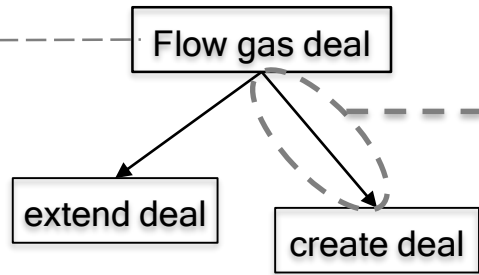
Main Entities



Behavioral Models

Activity Model

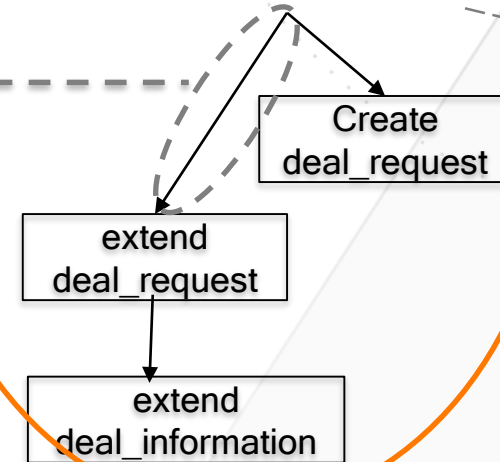
Activity event type
i.e. execution of one
activity included in
the email



=> event timestamp ≠
email timestamp

Activity & SA Model

Flow gas_information



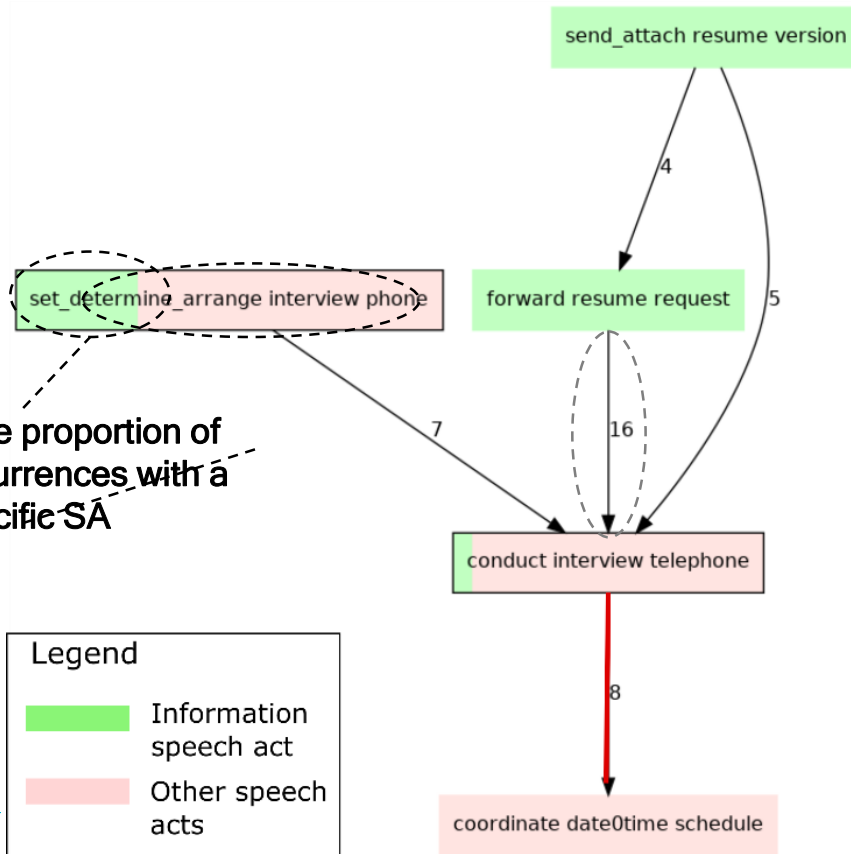
Activity & SA event type
i.e. sending an email containing
the activity expressed with a
specific speech act

=> event timestamp =
email timestamp

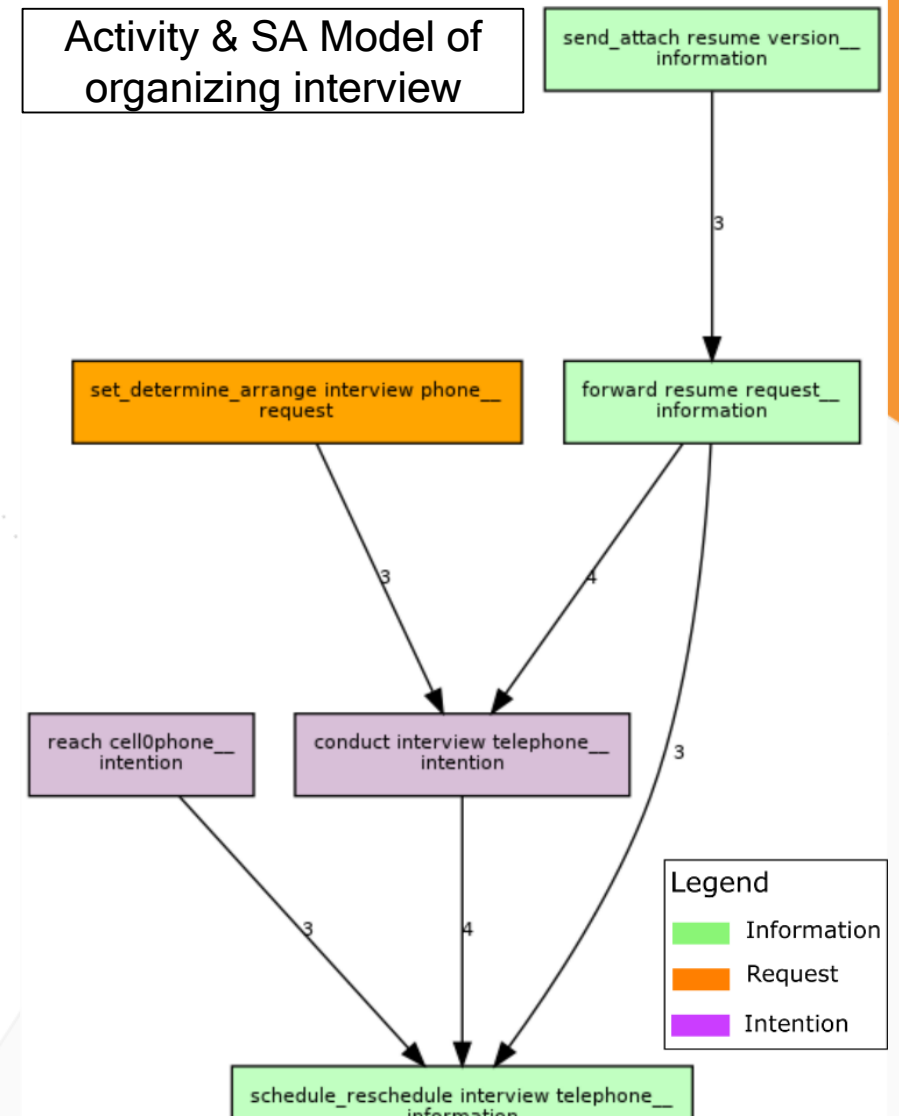
Better reflects how emails are
used in the context of performing
BP

BEHAVIORAL PERSPECTIVE: VISUALIZATION EXAMPLES

Activity Model of the BP fragment organizing interview



Activity & SA Model of organizing interview



EVENT LOG GENERATION: EVALUATION DATASET

- **Size:** 7056 emails from the public Enron dataset
- **Dataset composition:** Emails of employees having different organizational and business roles collected in the way that they form actor groups

Organizational roles of employees

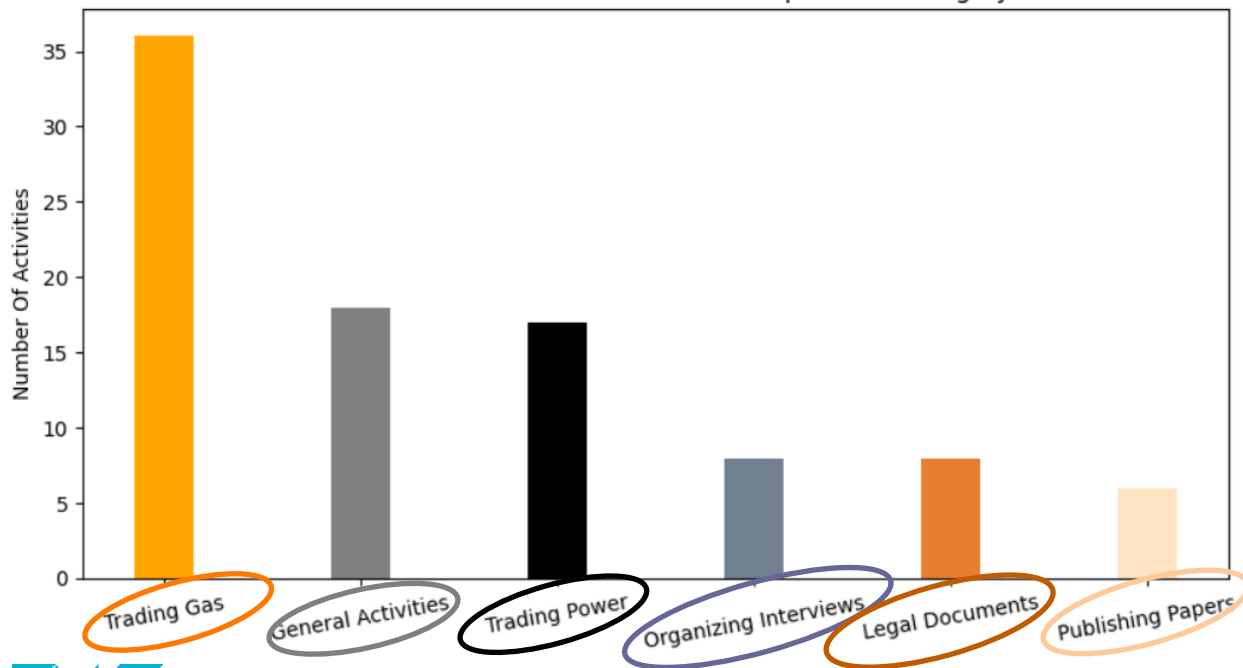
Business roles of employees

<i>Ep</i>	<i>OR</i>		<i>BR</i>	<i>N_{em}</i>	<i>N_{emRe}</i>	<i>N_{Re}</i>	<i>N_{emTh}</i>	<i>N_{reTh}</i>		
E1	Managing Director		Trading	343	421	18	1180	95		
E2	Senior Counsel		Legal	102						
E3	Managing Director		Risk	2283		8				
E4	Assistant		Management	357	682	31				
E5	Manager		Logistics	738						
E6	Specialist		Settlements	108						
E7	Specialist		Logistics	100						
E8	Employee		Employee	158						
E9	Specialist		Logistics	80						

EVENT LOG GENERATION: OVERVIEW ON THE DISCOVERED ACTIVITIES FROM THE EMAILS OF THE 9 EMPLOYEES

- Total number of discovered activities: 102
- % of Relevant activities: 93 %
- Example of activities (see the following table)
- Main categories:

Distribution of the number of activities per main category



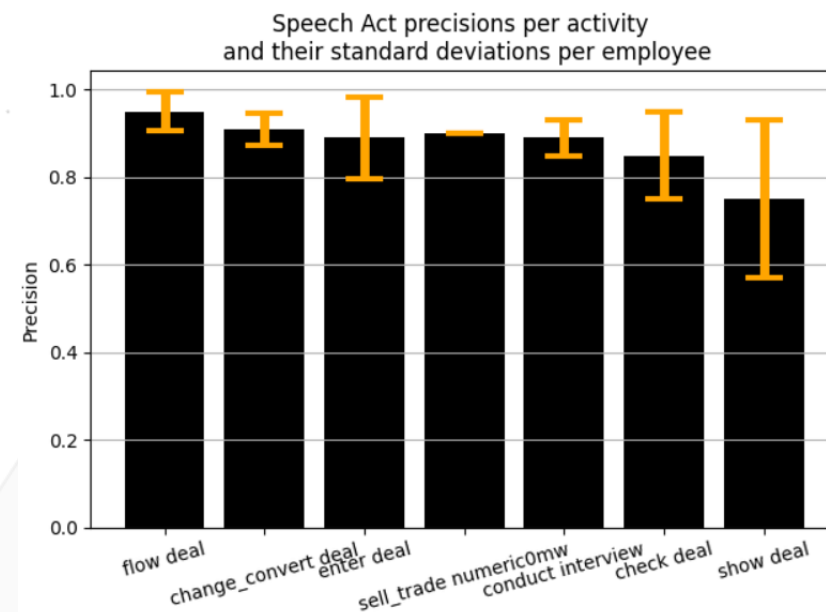
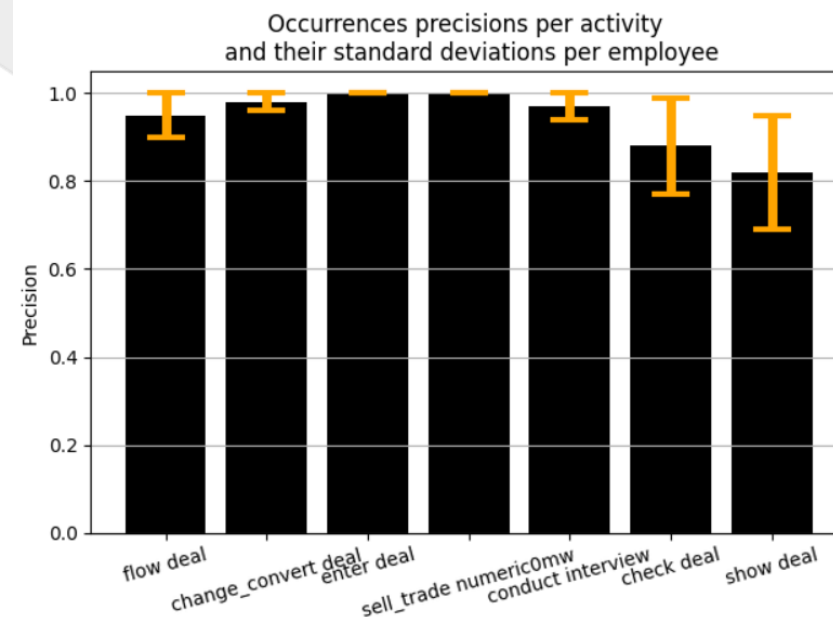
ID	Activity	N _{occ}
1	<i>flow deal{gas}</i>	189
4	<i>create deal{numeric}{ticket}</i>	113
5	<i>sell_trade numeric0mw{pricenumeric}</i>	99
7	<i>conduct interview{telephone}{informal}</i>	93
8	<i>extend deal{numeric}{rest}</i>	90
10	<i>send_attach resume{version}{electronic}</i>	81
12	<i>purchase_buy gas{plant}</i>	69
15	<i>set_determine_arrange interview{phone}</i>	60
20	<i>send_attach_ship transport0contract{deal}{term}</i>	49
24	<i>sell_trade natural0gas{plant}</i>	41
29	<i>purchase_buy numeric0mw{pricenumeric}</i>	34
30	<i>make reservation{hotel}</i>	33
32	<i>long numeric0mw{hour0end}</i>	27
35	<i>schedule_reschedule meeting{assistant}</i>	26
36	<i>send_attach info_information</i>	26
43	<i>attend meeting</i>	20
48	<i>deliver gas</i>	18
58	<i>write_publish book</i>	16
62	<i>short numeric0mw{ho0numeric}</i>	14
79	<i>execute agreement</i>	10
91	<i>register conference</i>	9
95	<i>review comment</i>	8

EVENT LOG GENERATION: EXPERIMENTS FOR EVALUATION

Experiment 2: Study the Overall Features of the generated event log

BP element	Number	Metric & Value
Activities	102	Relevance = 0.93
Activity occurrences	3102	Precision = 0.85
Speech acts	--	Precision = 0.88
Threads	1287	Consistency = 0.85
Relevant information values	194	

- It reflects to which degree each obtained thread refers to one instance;
- It corresponds to the average of the inverse number of the annotated instances per each obtained thread;

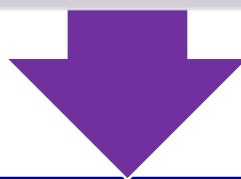


Research Problem

- How to enable cognitive process analysis on email data?

1. How to extract process related information from emails?

activities, actors, business artifacts, etc.



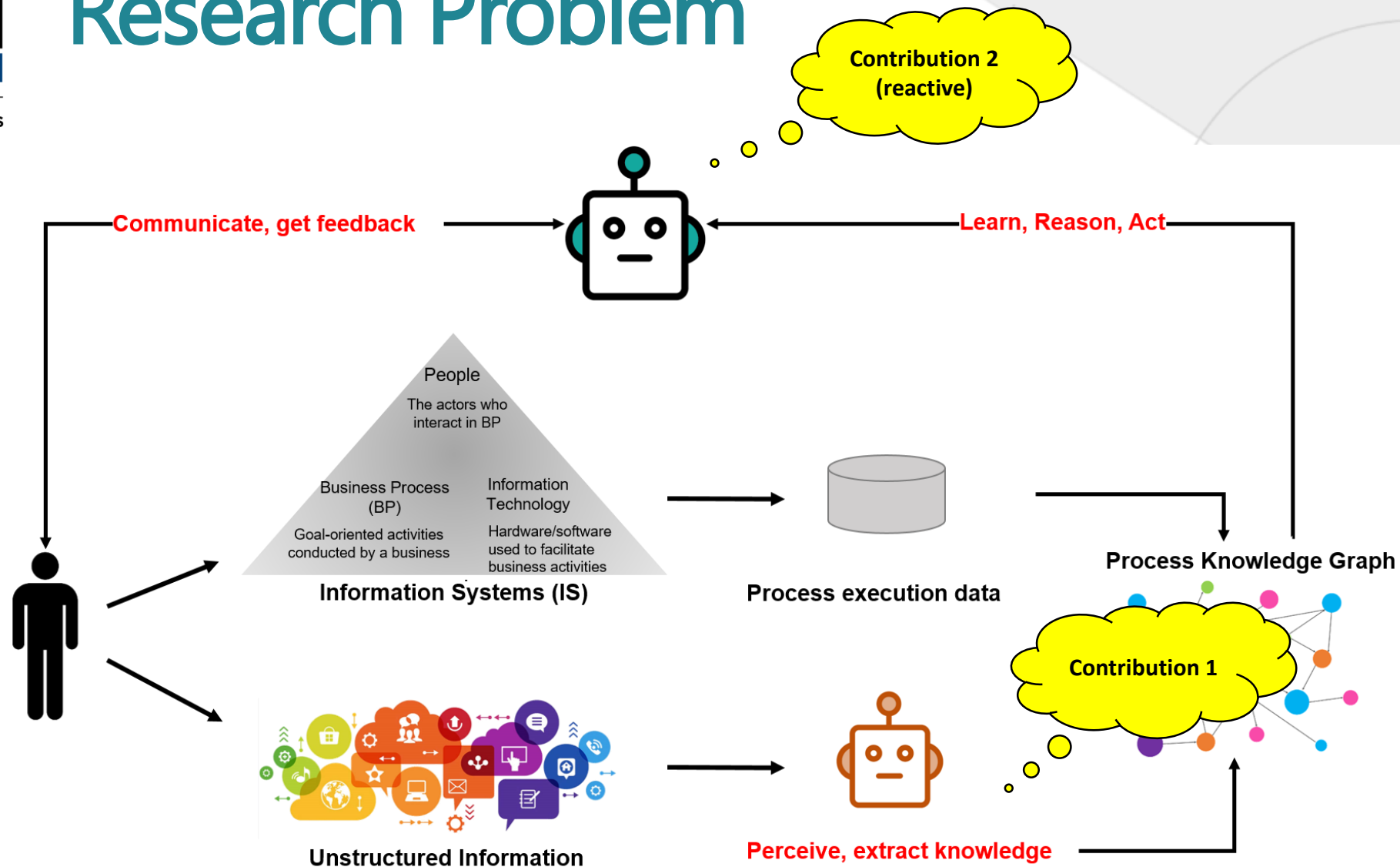
3. How to assist end users in querying the process data?

Querying the process data should be accessible to business experts

Contribution 1

Contribution 2

Research Problem



Research problem

- How to easily query the stored process data?
- Existing works can be categorized into:
 - **Rule-based approaches :**
 - +Pros:** Independent from training data + high generality and adaptability to new DB
 - Cons:** Has its own assumption about the query, NL question and DB schema.
 - **Machine-learning approaches :**
 - +Pros:** Allows the user to express his question in NL with great flexibility.
 - Cons:** This necessitates the use of training data, which is not always available.

We propose a hybrid NLI system that combines rule-based and machine learning approaches.

NL to cypher

- example

Applications with an amount greater than 10000 and validated by John

```
MATCH (actor: Actor)-[:HasContributed]-(activity: Activity),  
        (activity)-[:AffectArtifact]-(application: Application)
```

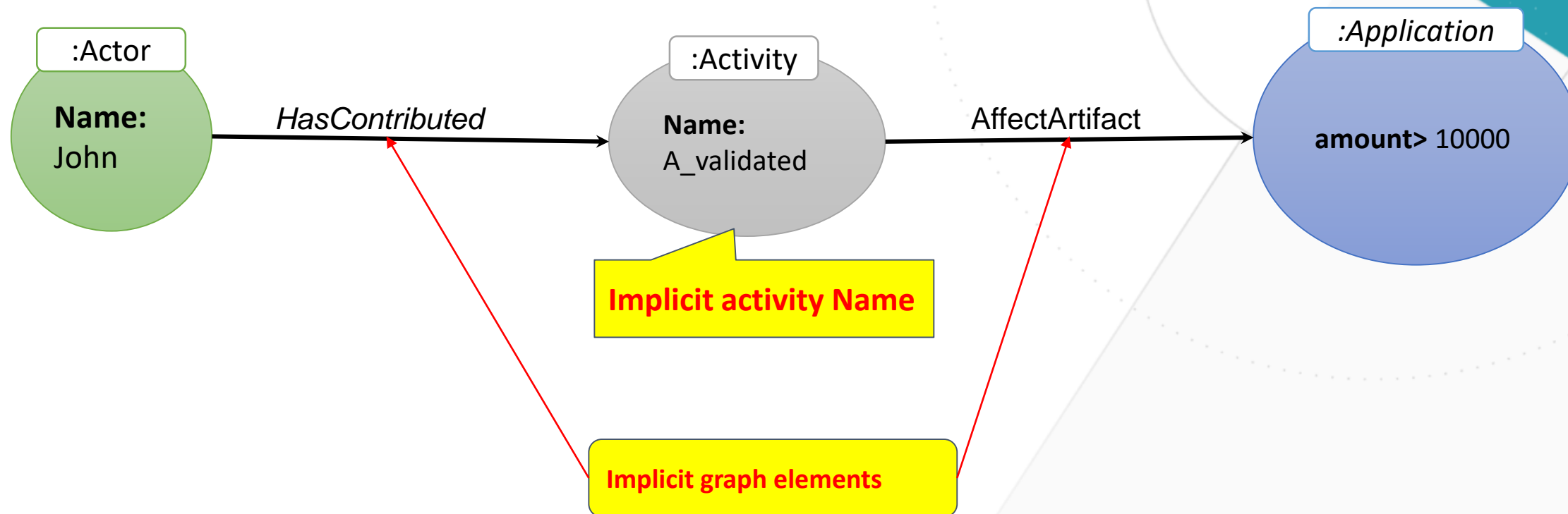
```
WHERE actor.Name= 'JOHN' AND application.amount> 10000  
        AND activity.Name= 'A_validation'
```

```
RETURN (application)
```

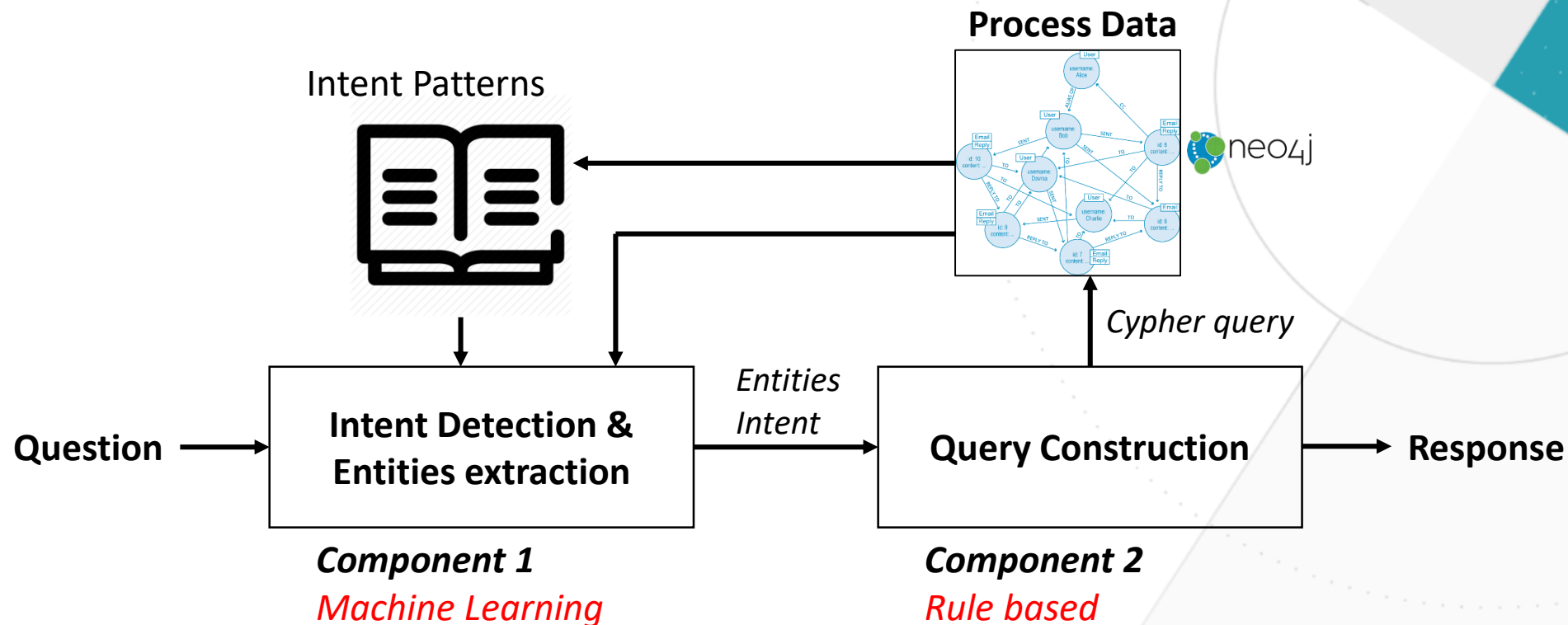
NL to cypher

- example

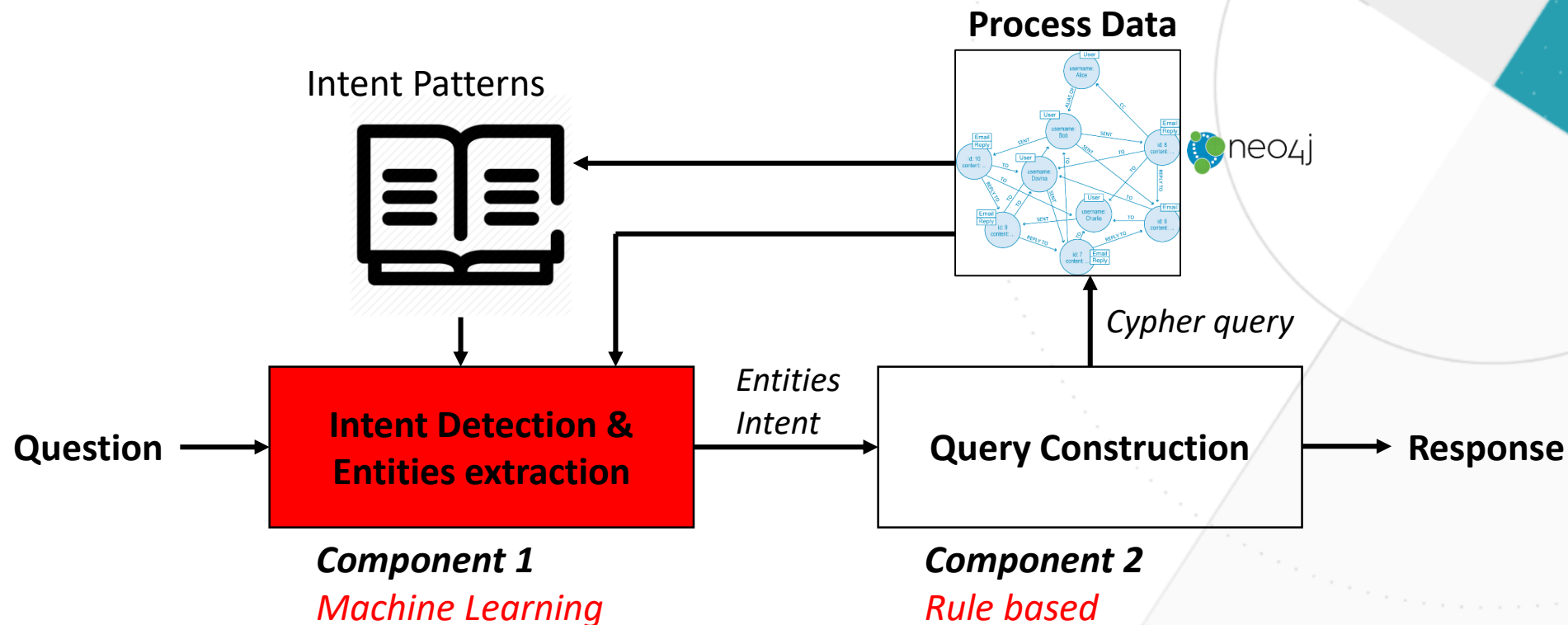
Applications with an amount greater than 10000 and validated by John



Bird's-eye view



Bird's-eye view



Component 1

- Intent detection

The intent describes the question purpose

- Minimal graph elements are required to construct the cypher query => **MATCH**
- What type of information should be returned => **RETURN**

We defined a set of **intent patterns** inferred from our meta-model

“Applications with an amount greater than 10000 and validated by John”

Application_AffectArtifact

“How many tasks were executed today?”

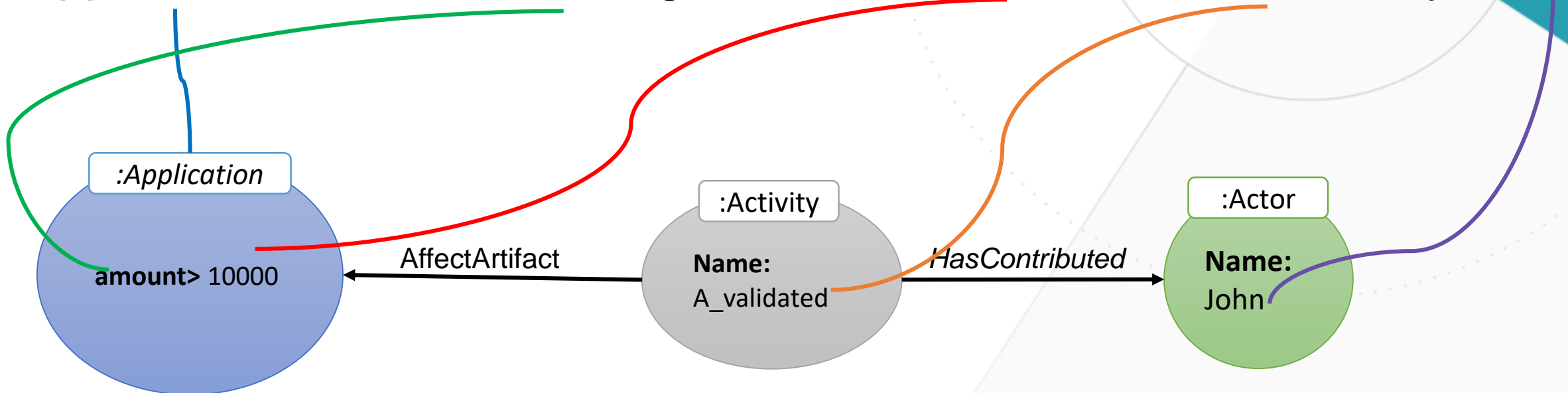
Activity_Count

Component 1

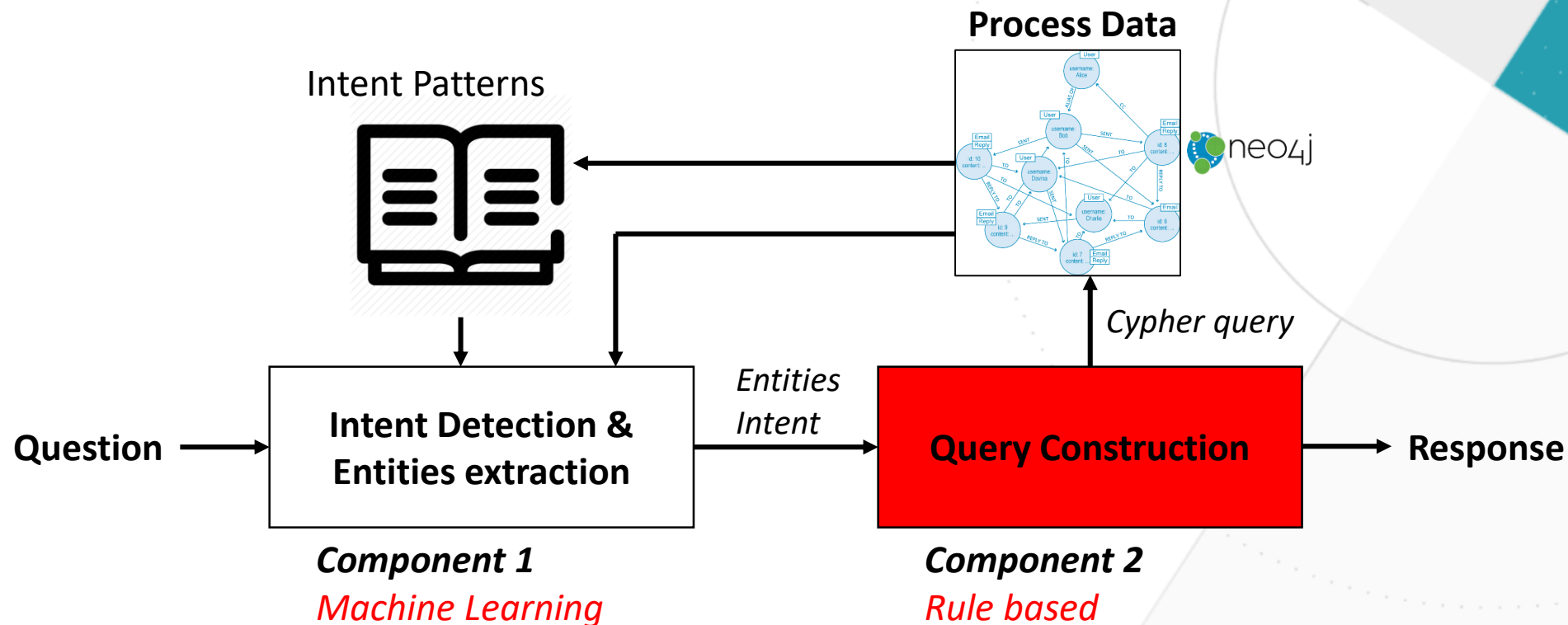
- Entities Extraction

- Entities are words that appear in the question and should be mapped to the graph elements (nodes, relationships, attributes, values).

Applications with an amount greater than 10000 and validated by John



Bird's-eye view



Component 2

- Query construction

Application type

Activity Name

Date

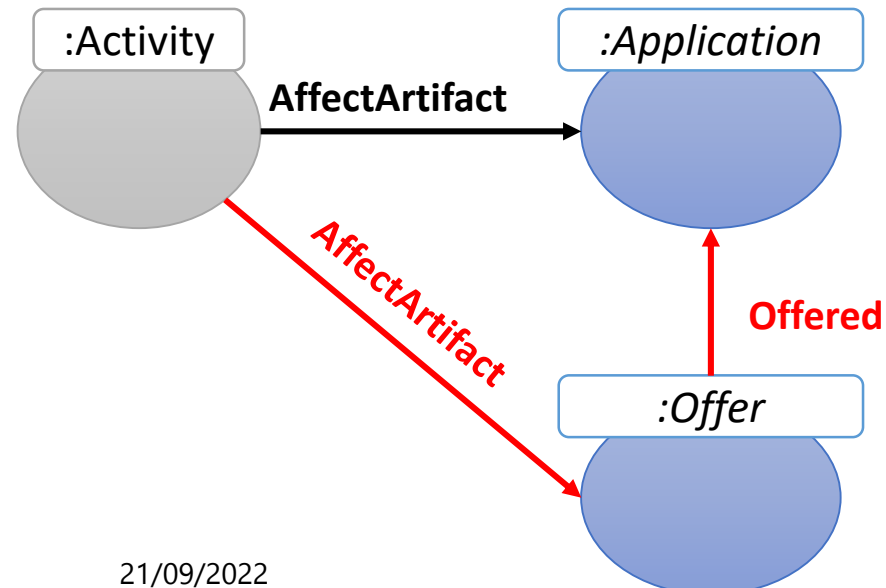
Return all **car** loan **applications** **validated** within the **last three months**, and which received **3 offers** minimum, order them by their **requested amounts**.

Number

Offer node

Application amount

- **Intent:** Application_AffectArtifact



MATCH (activity: Activity)-[:AffectArtifact]-(application: Application),
 (activity)-[:AffectArtifact]-(offer: Offer),
 (offer)-[:Offered]-(application)

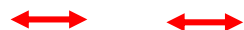
RETURN (application)

Component 2

- Query construction

WHERE application.Type= 'Car' **AND** activity.Name='A_validated' **AND** activity.Time> 14-10-2021

Return all **car** loan **applications validated** within the **last three months**, and which received **3 offers** minimum, order them by their **requested amounts**.



trigger word of greater than operator

WITH application, COUNT(offer) as offerCT

WHERE offerCT >=3

Component 2

- Query construction

Return all **car** loan **applications** **validated** within the **last three months**, and which received **3 offers** minimum, order them by their **requested amounts**.

trigger word of ORDER BY

Application amount

ORDER BY application.amount

Component 2

- Query construction

Return all **car** loan **applications validated** within the **last three months**, and which received **3 offers** minimum, order them by their **requested amounts**.

```
MATCH (activity: Activity)-[:AffectArtifact]-(application: Application),  
         (activity)-[:AffectArtifact]-(offer: Offer),  
         (offer)-[:Offered]-(application)
```

```
WHERE application.Type= 'Car' AND activity.Name='A_validated' AND activity.Time> 14-10-2021
```

```
WITH application, COUNT(offer) as offerCT
```

```
WHERE offerCT >= 3
```

```
RETURN (application)
```

```
ORDER BY application.amount
```

Achievements : contribution 1

- 1) Formalize the definition of BP knowledge that could be discovered from emails
- 2) Introduce a **totally unsupervised** approach for BP fragment discovery w.r.t multiple perspectives
 - **Without requiring priori information** concerning BP knowledge in emails => **minimize human intervention**
 - Composed of several algorithmic solutions for **event log generation** & **event log mining**:
 - Introduce a learning solution for **activity discovery** from emails based on discovering **low dispersed patterns of concepts** & grouping patterns **without requiring** the prior definition of the **number of activities**
 - Rely on **overlapping grouping of activities** to discover artifacts & BP fragments (Data and Function perspectives)
 - Discover **actor contributions** when performing activities (Organizational perspective)
 - Estimate the **event sequencing** constraints **in the absence** of precise information concerning **event timestamps** (Behavioral Perspective)
 - Validated using the **public** dataset Enron while **sharing the obtained results** to ensure a type of comparison with existing works

Achievements : contribution 2

- We proposed an **intent-based NLI for querying process** execution data.
 - facilitates the querying activity by understanding and interpreting the **intent of the user** from a natural language question,
 - constructs automatically the corresponding Cypher query to be executed over the process data stored in a **graph database**, and returning the answer
 - Validated using the **public** a real-life event log from BPIC'17.

Conclusion

- Cognitive Process Analytics system that **understands, reasons, discovers** actionable insights and **interacts**
- Where we are today

Understand



Operationalize structured/unstructured data, feed into process knowledge graphs

Reason



Reason over knowledge graphs (recent efforts by the community to enable process analytics on knowledge graphs)

Discover

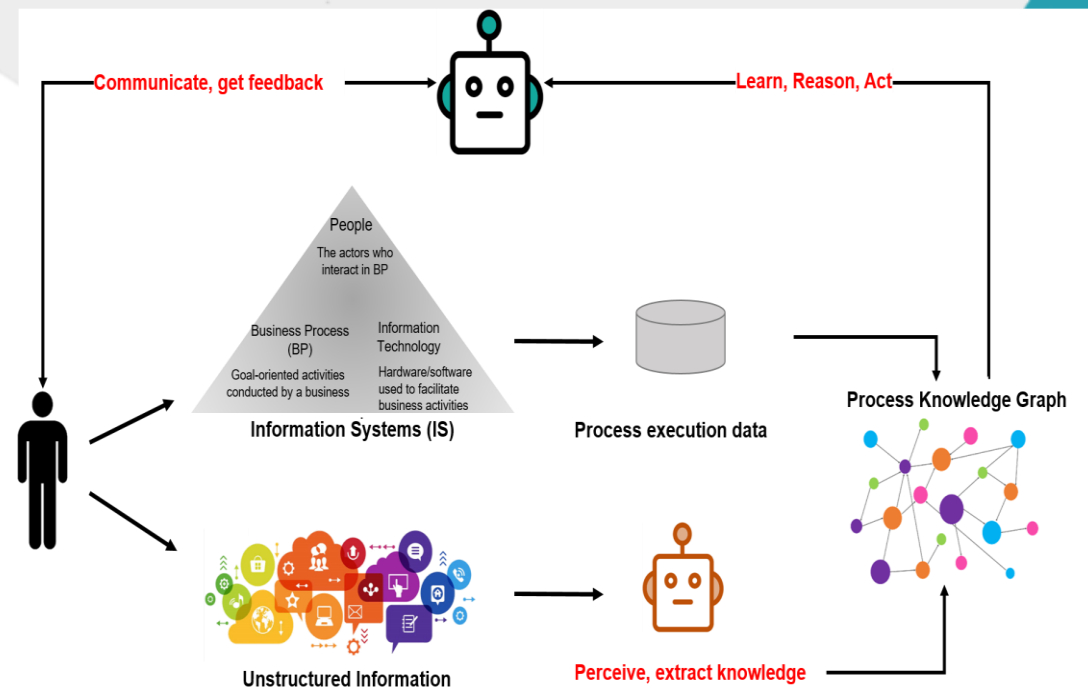


Reactive and proactive bots

Interact

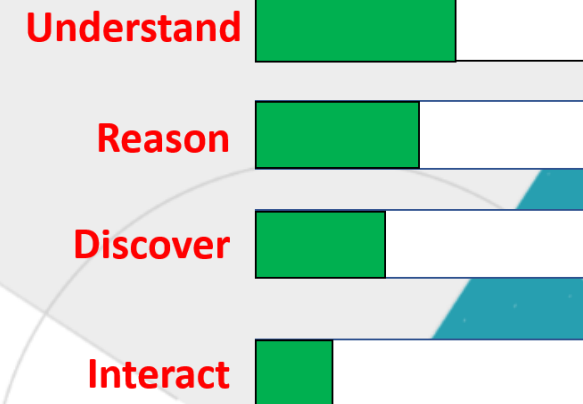


See, talk and hear (we are still in the beginning)



Research directions

- Enrich process knowledge graphs
 - Currently we capture information at the instance level
 - Add and connect extra layers (process, system, organization, etc.)
- Develop reasoning approaches over enriched knowledge graphs
- From reactive to proactive conversational bots
 - Current work on analyzing received emails and generating response templates



Email answering from process perspective

- An RPA approach to automate the generation of email replies when performing processes within emailing systems.

Goal:

- help employees access process oriented knowledge included in emails
- assist them in performing their repetitive process activities.

The provided recommendations are mainly:

1. Email responses templates,
2. The related process knowledge:
 - activities of emails responses,
 - speech acts,
 - the related business data values

References

- [1] M.E et al., **Discovering Activities from Emails Based on Pattern Discovery Approach.** [BPM \(Forum\) 2020](#): 88-104
- [2] M.E et al., **Discovery of Activities' Actor Perspective from Emails based on Speech Acts Detection.** [ICPM 2020](#): 73-80
- [3] M.E et al., **A Meta Model for Mining Processes from Email Data.** [SCC 2020](#): 152-161
- [4] M.E et al., **Discovering Business Processes And Activities From Messaging Systems: State-Of-The Art.** [WETICE 2020](#): 137-142
- [5] **Multi-perspective business process discovery from messaging systems: State-of-the art.** [CCPE Journal](#), 2021, p. e6642.
- [6] N.L, M.E et al., **Emails Analysis for Business Process Discovery.** [ATAED@Petri Nets/ACSD 2019](#): 54-70
- [7] M. K, et al.: **An Intent-Based Natural Language Interface for Querying Process Execution Data.** [ICPM 2021](#): 152-159

THANK YOU