

## KDD<sup>1</sup> – Overcoming Massive Data Streams for Intelligence Tasks

**Dr. Vera Kamp**  
PLATH GmbH  
Gotenstrasse 18  
20097 Hamburg  
Germany

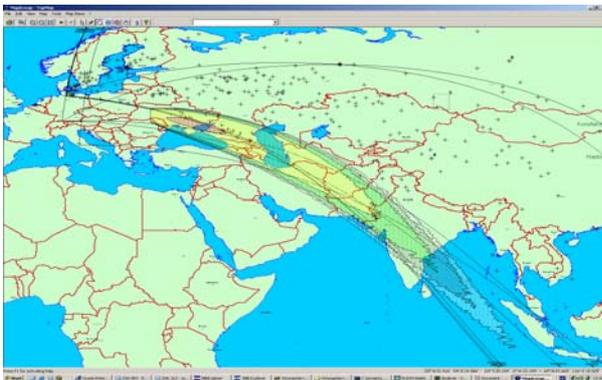
[vera.kamp@plath.de](mailto:vera.kamp@plath.de)

### **ABSTRACT**

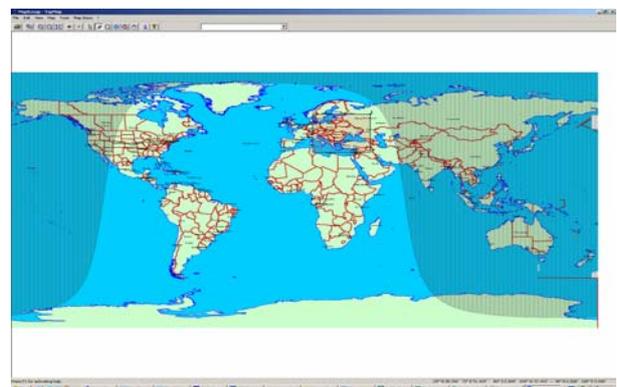
*Actually very interesting IT systems promise to reveal connections between apparently harmless and unrelated information pieces. An article from the New York Times in February 2006<sup>2</sup> makes clear that common data mining techniques were not successful in general. Despite huge investments, correlating data from different sources did not yield satisfactory results. Transforming low-level data by aggregation to meaningful events is nevertheless the key to building the basis for succeeding decisions in the context of situation reports*

*More realistic and manageable is an approach that includes interactions with the user along with domain specific knowledge. Gaining security relevant messages should be based on an iterative multi-level process. This process represents the core element of intelligence analysis systems which play an important role for supporting decisions in management information systems<sup>3</sup>.*

*The following example illustrates the principal automated process for discovering communication structures in the context of radio reconnaissance: A crucial part of this process is the analysis and visualisation of communication structures, or more generally, of network information. This should be embedded in spatio-temporal data analysis with geo-oriented data access and the integration of domain-specific analysis functions.*



**Figure 1: Domain-specific Analysis Function**



**Figure 2: Spatial Access**

<sup>1</sup> Knowledge Discovery in Databases

<sup>2</sup> J. Markoff, The New York Times, nytimes.com, 25.2.2006: "Taking Spying to Higher Level, Agencies Look for More Ways to Mine Data"

<sup>3</sup> Jan Herring, "What is Intelligence Analysis?", Competitive Intelligence Magazine, Vol. 1. No.2, July-Sep., 1998, p.14 .

# Report Documentation Page

Form Approved  
OMB No. 0704-0188

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

1. REPORT DATE <b>01 DEC 2006</b>		2. REPORT TYPE <b>N/A</b>		3. DATES COVERED <b>-</b>	
4. TITLE AND SUBTITLE <b>KDD Overcoming Massive Data Streams for Intelligence Tasks</b>				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>PLATH GmbH Gotenstrasse 18 20097 Hamburg Germany</b>				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release, distribution unlimited</b>					
13. SUPPLEMENTARY NOTES <b>See also ADM002067., The original document contains color images.</b>					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT <b>UU</b>	18. NUMBER OF PAGES <b>32</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>			

The intelligent analysis of radio emission data is based on data mining techniques, cluster visualisations to validate the results, a model based communication detection (including domain-specific knowledge) and the visualisation of communications. The following use case of a simple simplex communication clarifies the problems and the applied methods. Module coupling is realised by a distributed architecture. Given are a huge amount of radio emissions which are arbitrarily distributed. Each emission is described by the attributes ID, frequency, modulation type, starting time, end time, latitude and longitude. It has to be considered that the data quality of single emissions depends on propagation conditions. Because these can vary, it can happen that single emissions or attributes are missing or on the other hand different classification level information are available. Furthermore, with a broadband collection of emissions the amount of information is extremely large and requires massive data handling which can not be processed in main memory.

**1.0 USE-CASE SIMPLEX-COMMUNICATION**

The use case is looking for a simplex communication chain with two stationary partners – a central station and a substation. Both are using the same constant nominal frequency and the same transmission mode. The partners are communicating alternating one after the other. The problem lies in the amount of possible communication structure instances. Although the communication can be easily described in an informal way it is necessary to find an exact, formal specification in order to perform a computer-supported analysis. It should not be realised by a specific static algorithm but should be interactively and exploratively changeable by the user. The core concept includes the following steps:

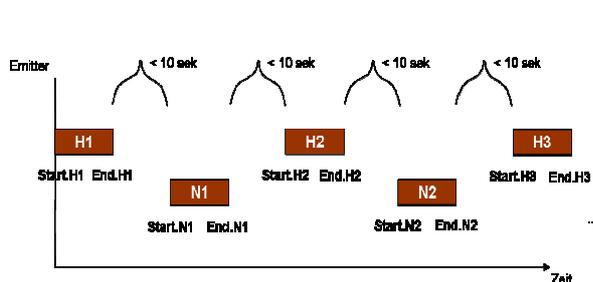
**1.1 Data Mining**

During the first step emissions are assigned to clusters. These subsume emissions concerning the spatial, temporal or frequency criteria. In this way significant data reduction is achieved. By spatial clustering special emitter station could be determined. Besides when processing of extremely huge data amounts the main problem to solve is how to choose the best method and parameters.

**1.2 Cluster Visualisation**

The next step serves the validation of the data mining results and already provides a possibility to manually discover communication structures by the user relying on the presented visualisation, for example the presentation of spatial clusters. Emission can appear as single instances or as temporal ordered parts of a cluster. It is difficult to visualise the emissions and clusters clearly arranged in order to focus on the actual interesting data. Additionally different attributes have to be integrated.

**1.3 Model based communication detection**



Computing communication structures from clusters is the next step. This is done by using typical communication models. A domain specific modelling language provides the possibility to represent the communication models. By this language the simplex communication can be formally specified. The model distinguishes between connection constitution and alternating communications.

**Figure 3: Simplex Communication Rules**

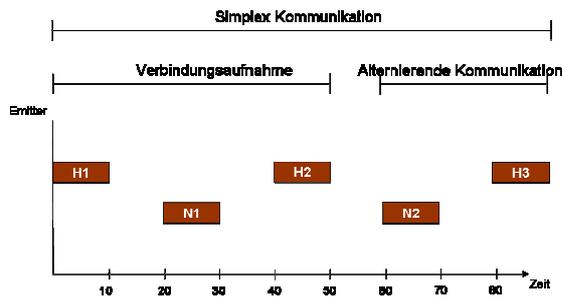


Figure 4: Simplex Communication Instance

The connection constitution consists of three emissions: the central station is sending, the substation replies. The alternating sequence consists of emissions of the central and the sub station. All emissions occur with the same frequency and modulation type. The distance between the emission is flexible by a delay parameter. A graphical notation of such a model is illustrated by the two adjoining pictures.

### 1.4 Visualisation of discovered communications

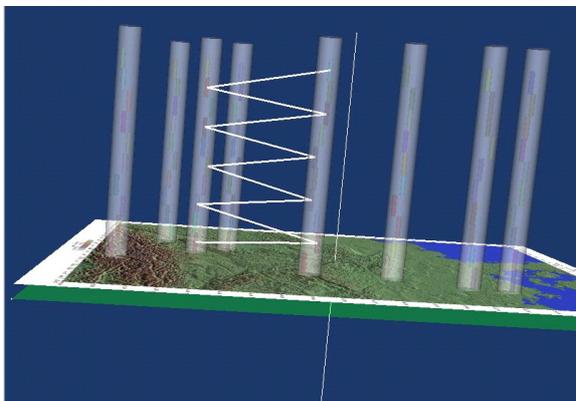


Figure 5: Discovered Network Communications

This step provides a presentation of the discovered communications and allows by this a validation of the model based communication detection. It has to deal with many composite events. A simple textual visualisation does not meet the needs. The graphical visualisation offers a better overview and manifold interaction possibilities.

The emitters are on the spatial level, time is the third dimension. The connection lines indicate the communications.

Overcoming massive data streams for intelligence tasks is a challenge which should involve the analysis process with a seamless data access and the intelligence analyst. The acceptance of the results depends on the possibility to validate the results. The sustainability of results has to be guaranteed by flexible extension of actual domain specific analysis methods.

## 2.0 REFERENCES

- [1] J. Markoff, The New York Times, nytimes.com, 25.2.2006: "Taking Spying to Higher Level, Agencies Look for More Ways to Mine Data"
- [2] Jan Herring, "What is Intelligence Analysis?", Competitive Intelligence Magazine, Vol. 1. No.2, July-Sep., 1998, p. 14.



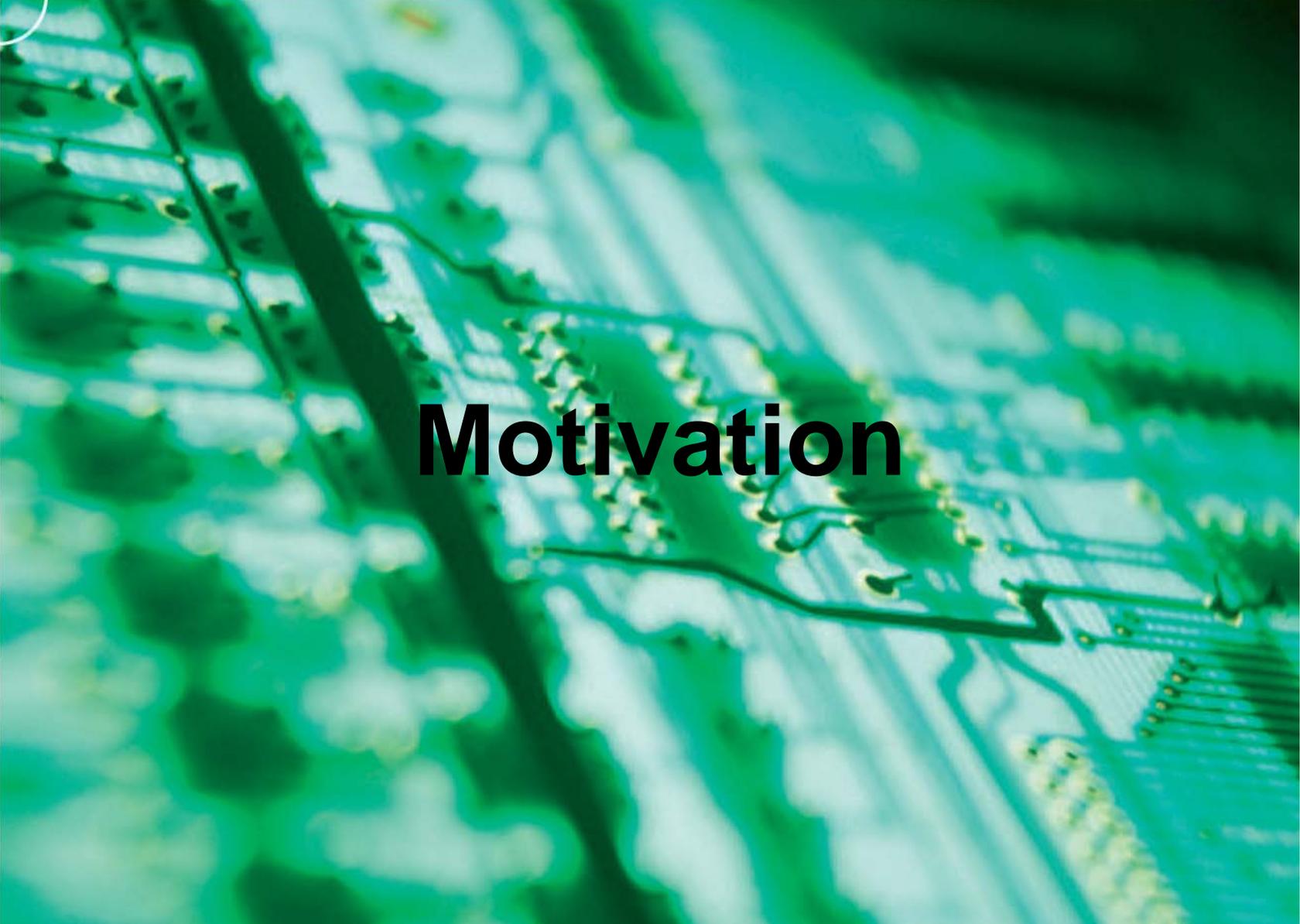


# KDD

## Overcoming massive data streams for intelligence tasks

Dr. Vera Kamp  
Dr. Joachim Stamm

- Motivation
- Basic Concept
- Exploration Process
- Summary



# Motivation

- ACOS delivers huge amount of data
- Intelligence requires to add value
- Task: Transforming low-level data by aggregation to meaningful events is the key to build the basis for succeeding decisions in the context of situation reports.
- Problem: “Common data mining techniques were not successful in general. Despite huge investments, correlating data from different sources did not yield satisfactory results.”

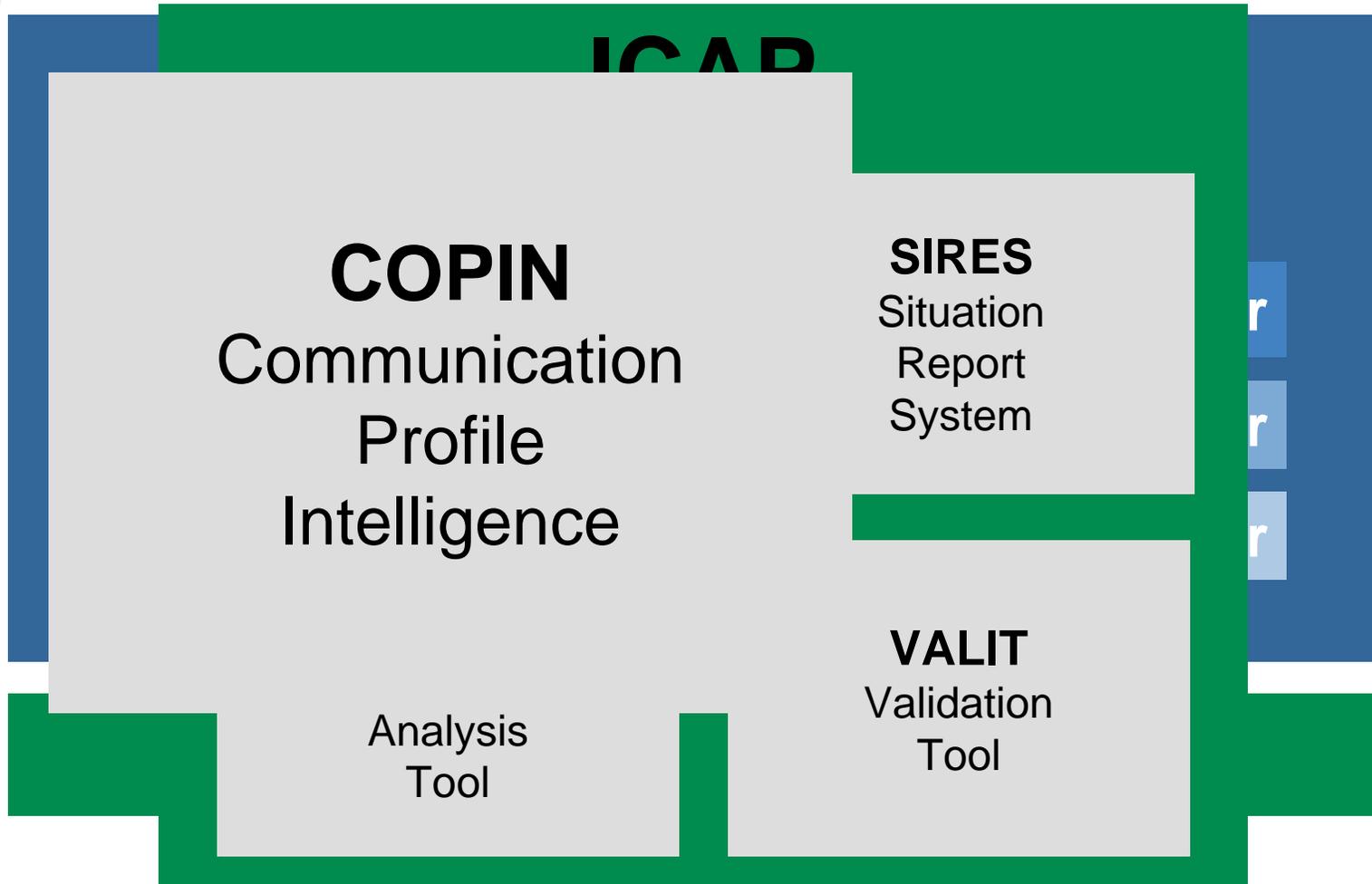
(New York Times in February 2006, Taking Spying to Higher Level, ...)

Domain specific approach is nessecary

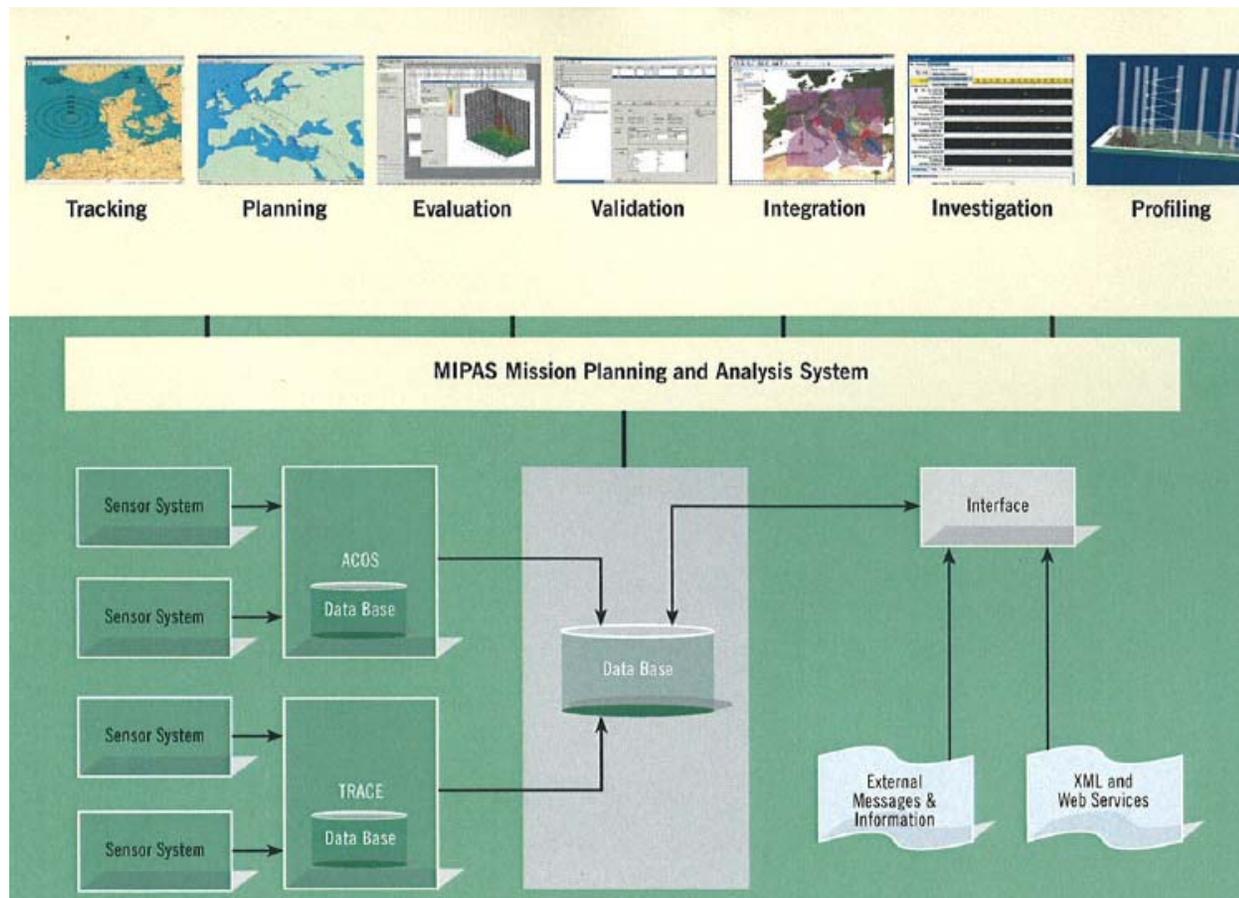
- Solution: **KDD - Overcoming massive data streams for intelligence tasks**
- Domain specific approach – is based on reconnaissance know-how
- More realistic and manageable: interactions with the user along with domain specific knowledge.
- Gaining security relevant messages should be based on an iterative multi-level process. This process represents the core element of intelligence analysis systems which play an important role for supporting decisions in management information systems.

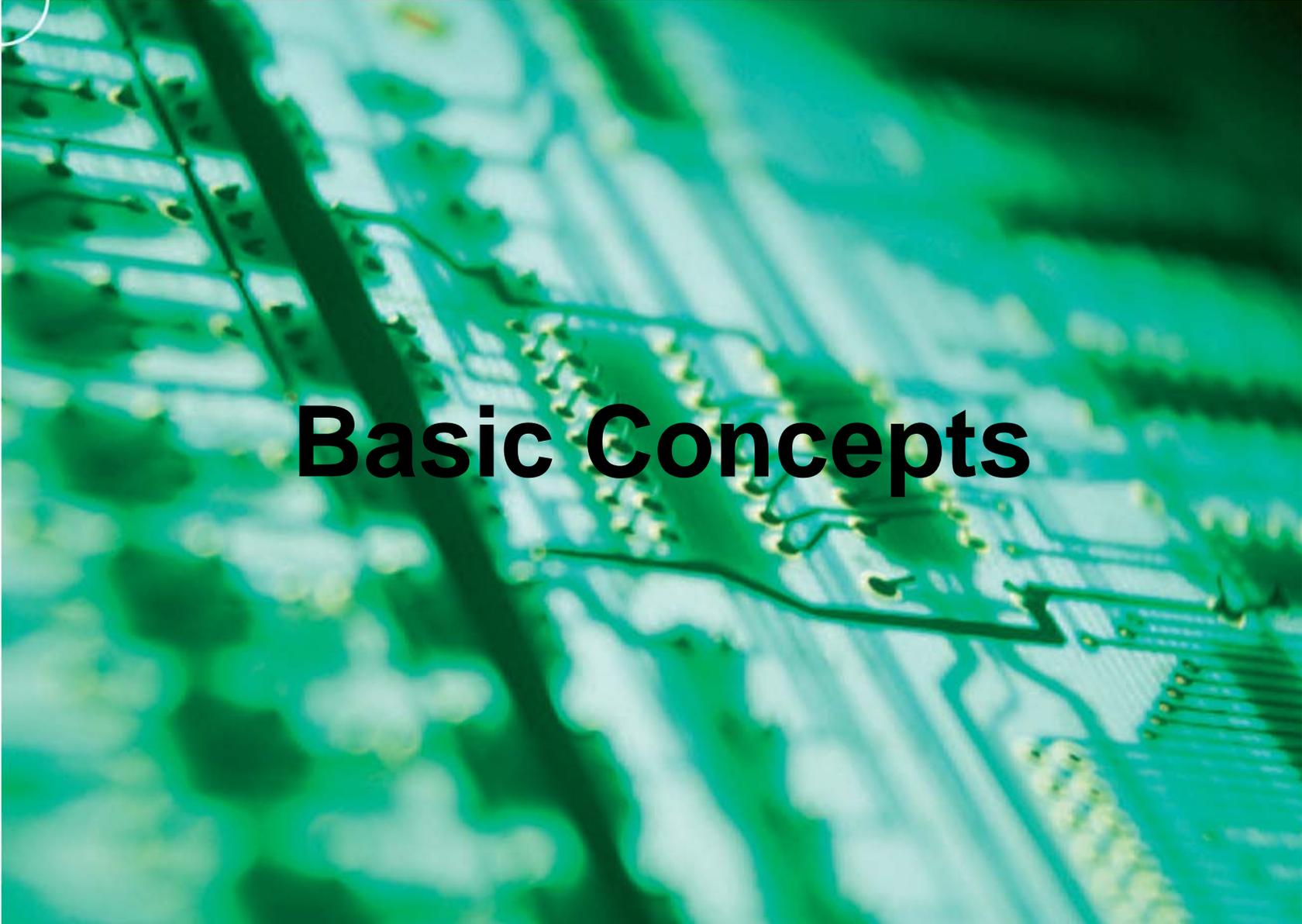
# MIPAS-Architecture Intelligent Evaluation Software

**PLATH**  
a clear signal



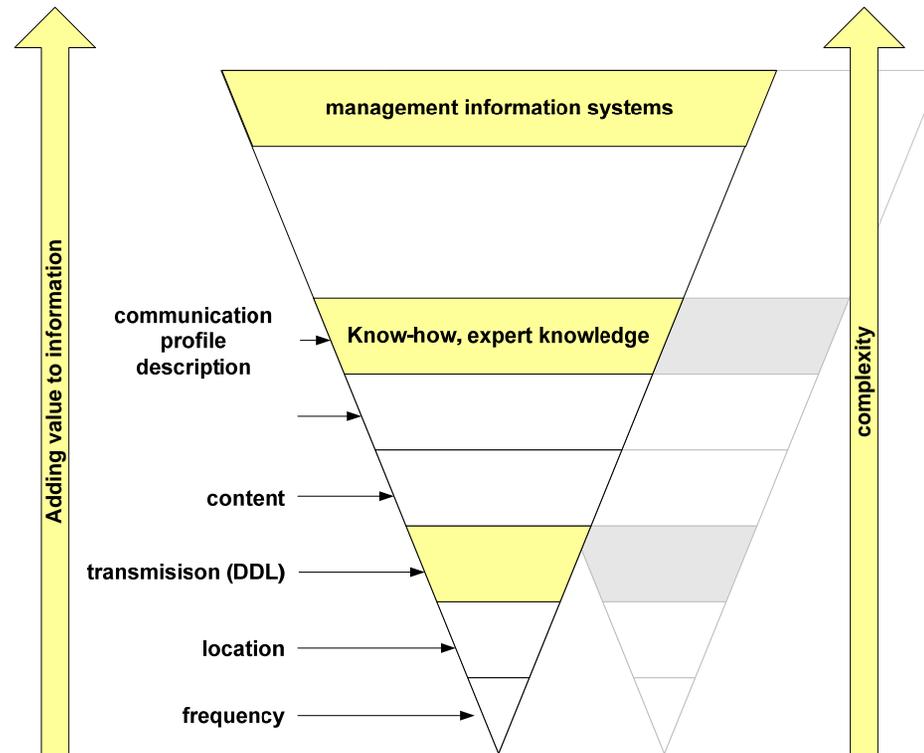
## C4ISR System for RF-Intelligence including Scalable & Modular Evaluation Software



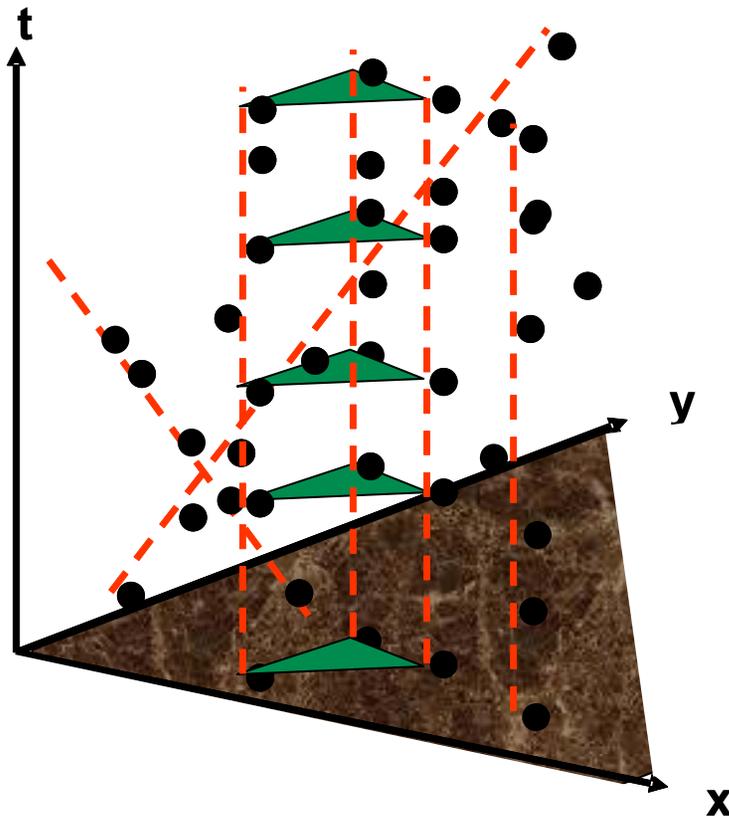


# Basic Concepts

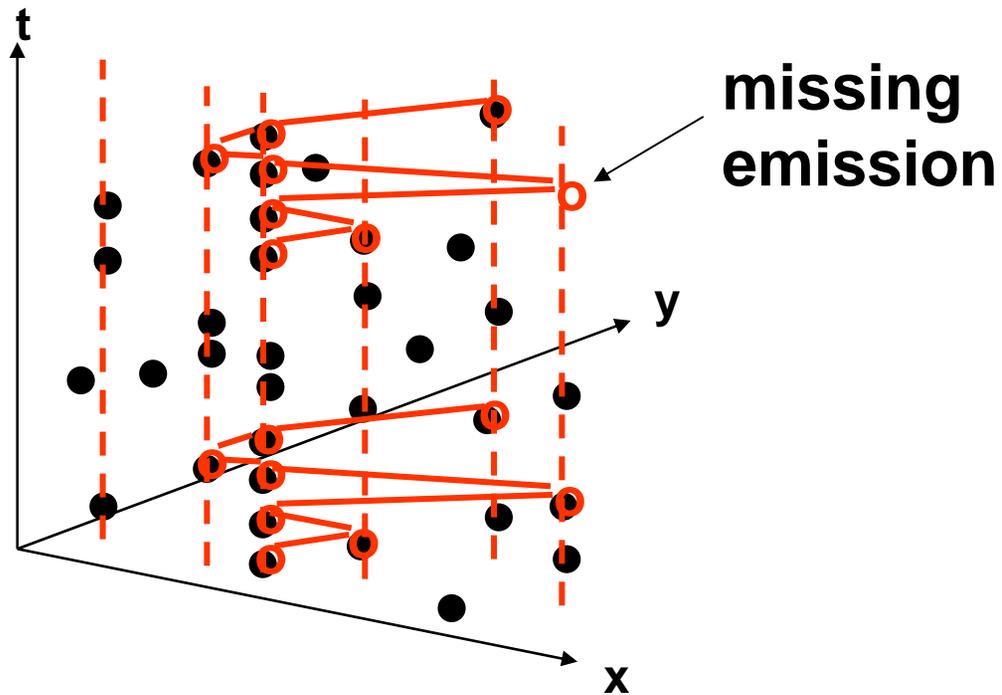
- ACOS provides for the 1st time:
  - broadband complete coverage
  - collection of all emissions on air are stored in database
- How can this information be exploited ?
- What are possible interests of our customers ?
- Does the data contain recurrent patterns ?
  - spatio-temporal
  - communication profiles
- Is there a deviation of normal behaviour
  - frequency change
  - new net members
- Is there an indication for the outfall of expected events ?
  - periodicity
  - member does not communicate



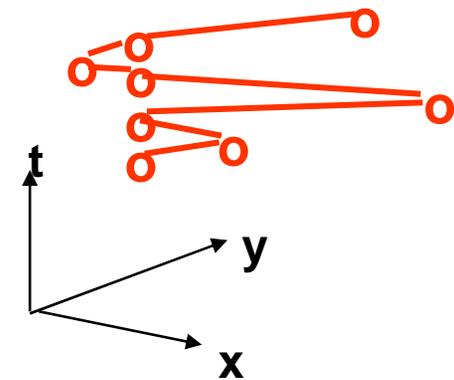
- Information processing aims at adding value to information
- Reduction to statements guides the development
- Important topics
  - Visualisation, Tracking/GIS
  - DM/KDD/OLAP
  - Data Warehouse

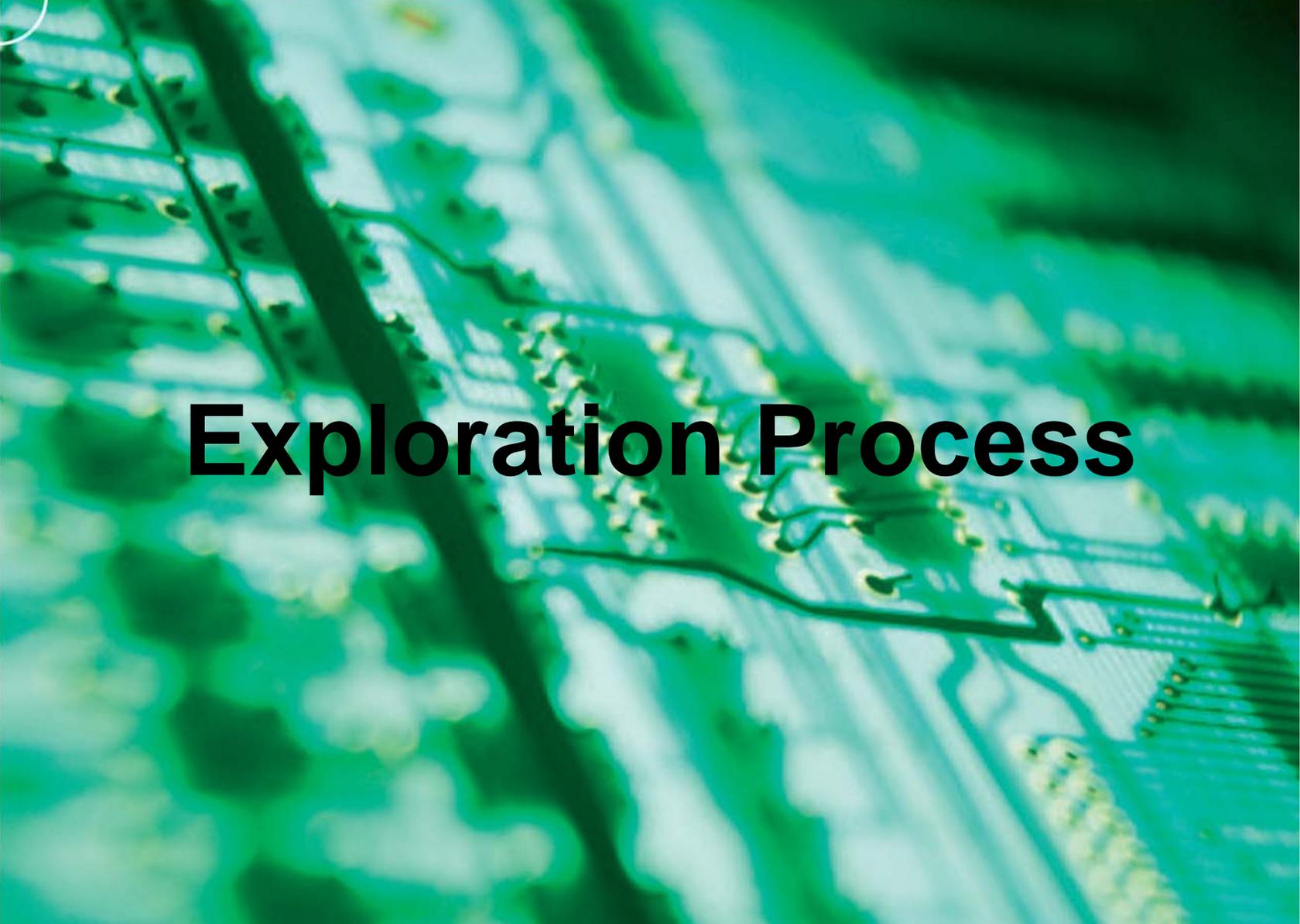


- **Discovery of stationary emitters**
- **Discovery of mobile emitters**
- **Discovery of simple command structures**



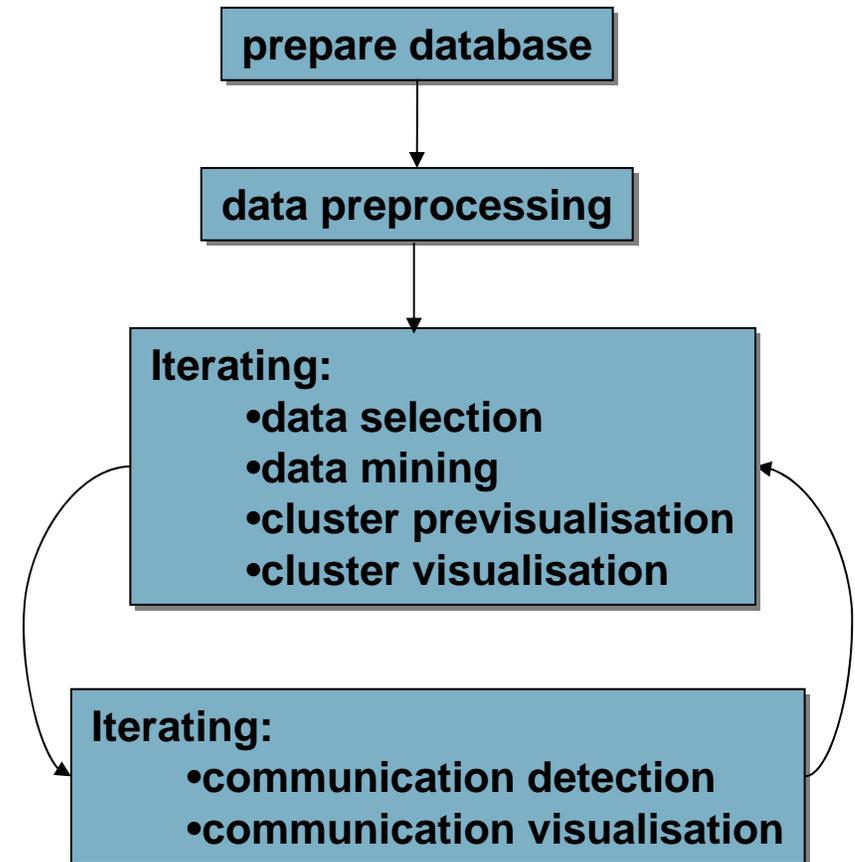
**model:**



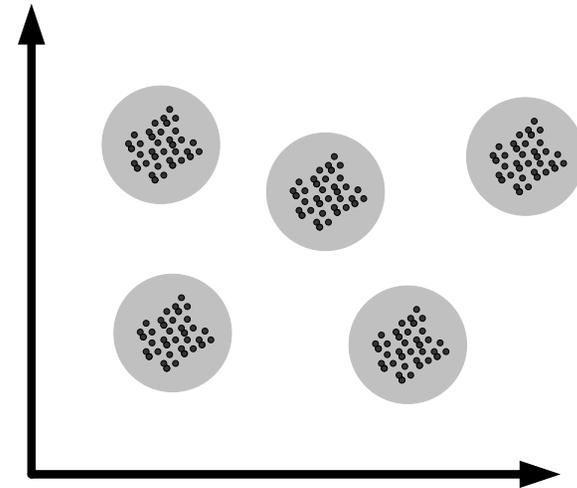


# Exploration Process

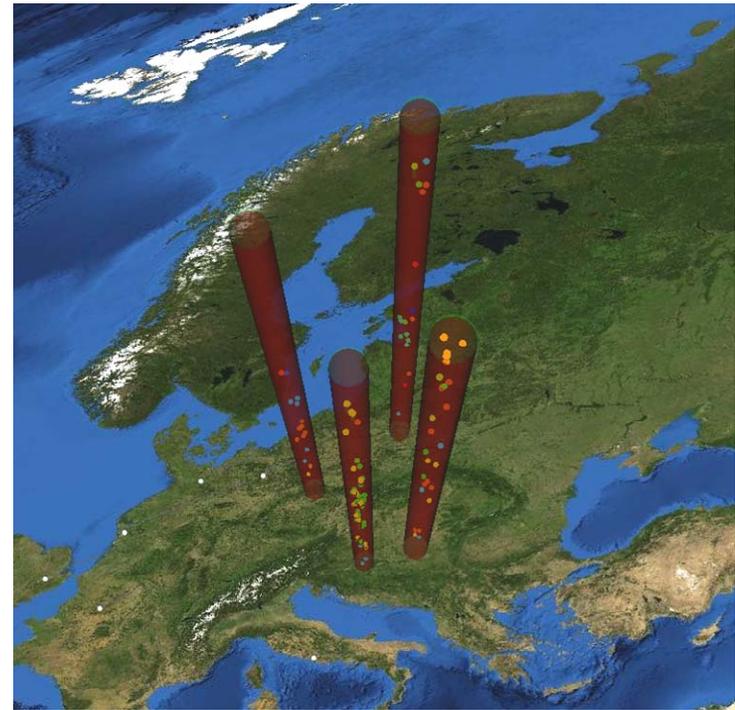
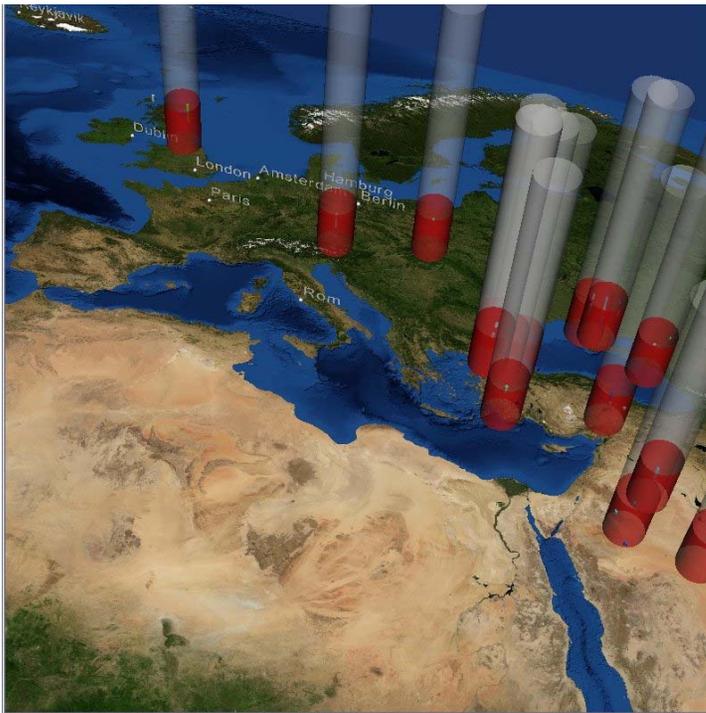
- 1. data mining**  
emissions → cluster  
(e.g. spatial combined emissions)
- 2. cluster visualisation**
- 3. model based detection**  
cluster → communication  
structures
- 4. visualisation of  
communication structures**



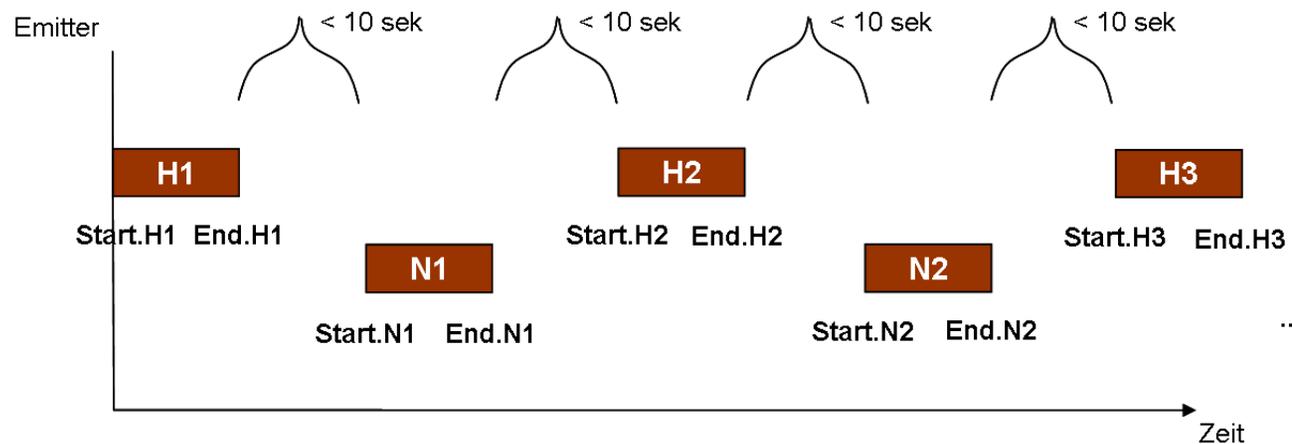
- clustering emission data:
  - location
  - frequency
  - time
- data reduction
- determination of emitters
  - spatial clustering
  - timeseries of emissions



# cluster visualisation



- modeling of typical communication structures
  - domain specific modeling language
  - representation of communication models
  - combination of model components
- event detection
  - translation of models in constraints
  - combinatorial analysis of possible model interpretation
- example



# Domain specific model based communication detection

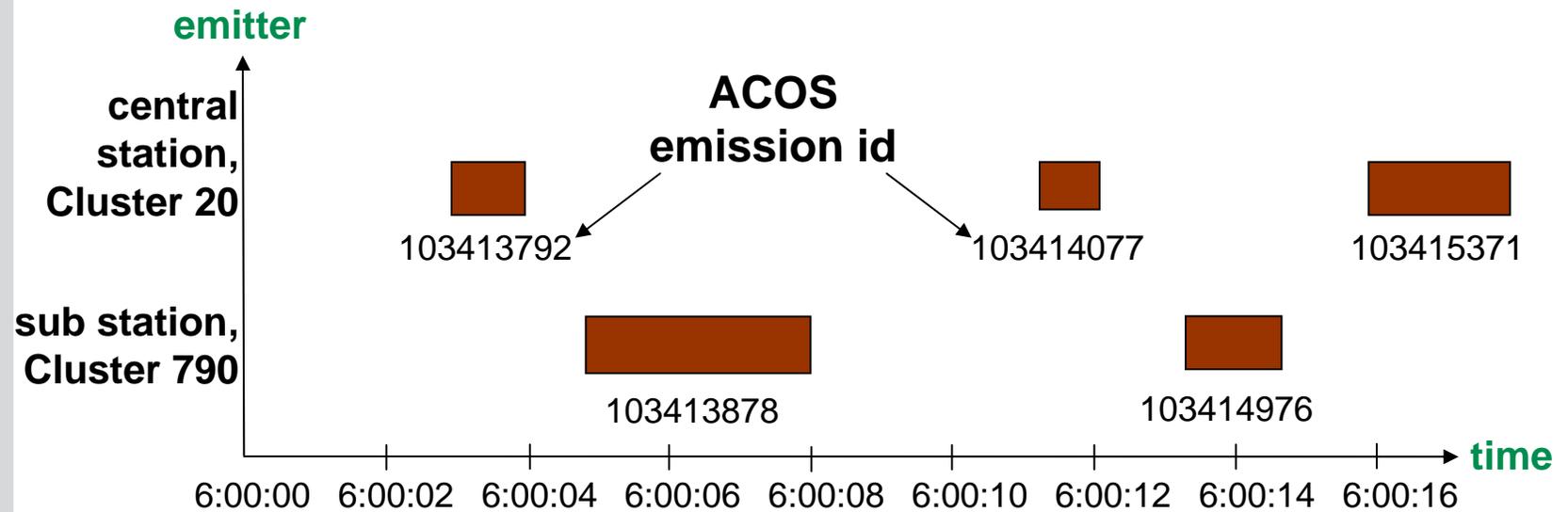
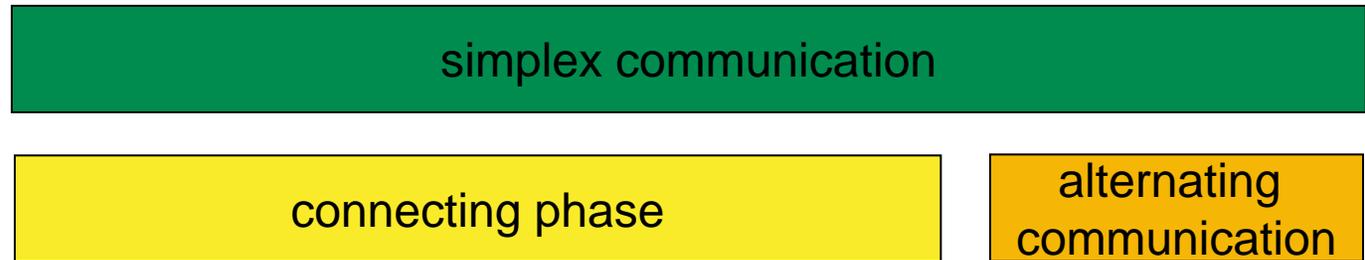
- central station: H
- sub station: N
- connection phase  
(3 emissions: H1, N1, H2)
- $H1.end < N1.start$ ;  
( $N1.start - H1.end$ )  $\leq 10$
- $N1.end < H2.start$ ;  
( $H2.start - N1.end$ )  $\leq 10$
- alternating communications  
(every 2 emissions: H3, N2)
- $N2.end < H3.start$ ;  
( $H3.start - N2.end$ )  $\leq 10$

The screenshot displays the 'Model Editor' software interface. At the top, there is a menu bar with 'File', 'Partner', and 'Communication'. Below the menu, there are checkboxes for 'Free Communication', 'Alternating Communication', and 'Consecutive Communication'. A yellow timeline at the top shows time intervals from 0 to 190 seconds, with markers every 10 seconds. The main area is divided into several sections:

- Headquarter:** Shows a frequency range of 0-1000 KHz, a 10 KHz bandwidth, and a location of long:135°/lat:48°. It includes a grid of 10 time slots, each with a duration of < 15s.
- Communication Partner 2:** Shows a frequency range of 500-1000 KHz, a 10 KHz bandwidth, and a location field. It includes a grid of 10 time slots, each with a duration of < 15s.
- Communication Partner 3:** Shows a frequency range of 10 KHz and a location field. It includes a grid of 10 time slots, each with a duration of < 15s.
- Communication Partner 4:** Shows a frequency range of 10 KHz and a location field. It includes a grid of 10 time slots, each with a duration of < 15s.
- Communication Partner 5:** Shows a frequency range of 10 KHz and a location field. It includes a grid of 10 time slots, each with a duration of < 15s.

Below the timeline, there are tabs for 'Frequency', 'Time', and 'Location'. The 'Partner Information' section includes a dropdown menu for 'Select Partner' (set to 'Communication Partner 2'), a 'Min. Frequency' slider (set to 500) and input field, and a 'Max. Frequency' slider (set to 1000) and input field. The 'Communication Information' section includes an 'Uncertainty' input field (set to 10). An 'Apply' button is located at the bottom right. A 'Help Information' section at the bottom provides instructions on how to use the frequency sliders and input fields.

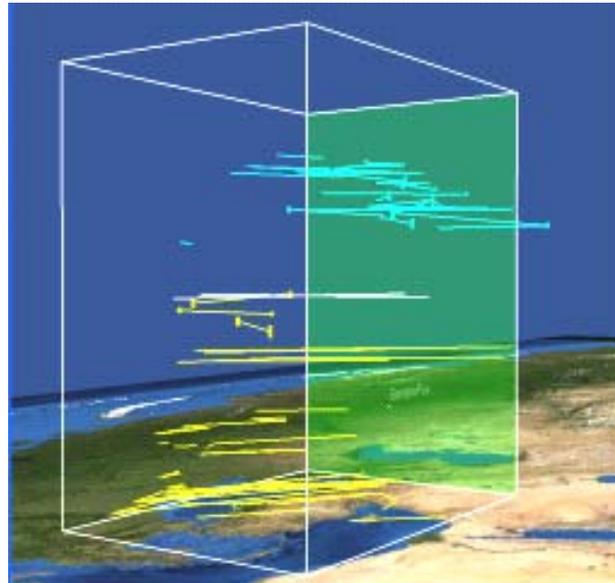
# model based event detection



frequency: 9991000.0, ModType: MFSK

# visualisation of a communication structure

## example: frequency change



We have:

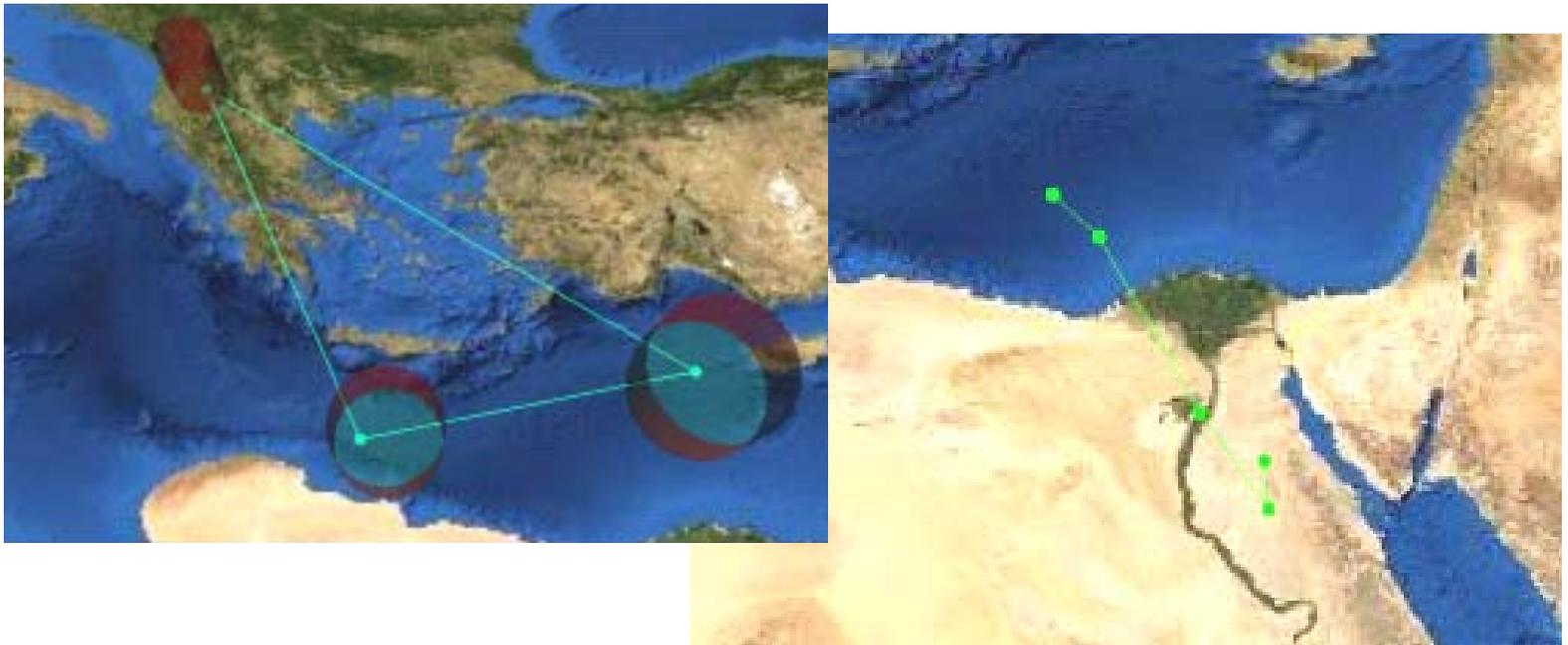
communications in an interesting time and space window with possibly known frequency ranges

We are looking for:

changes in frequency usage for the interesting communications

# visualisation of a communication structure

## example: command structures



We have:

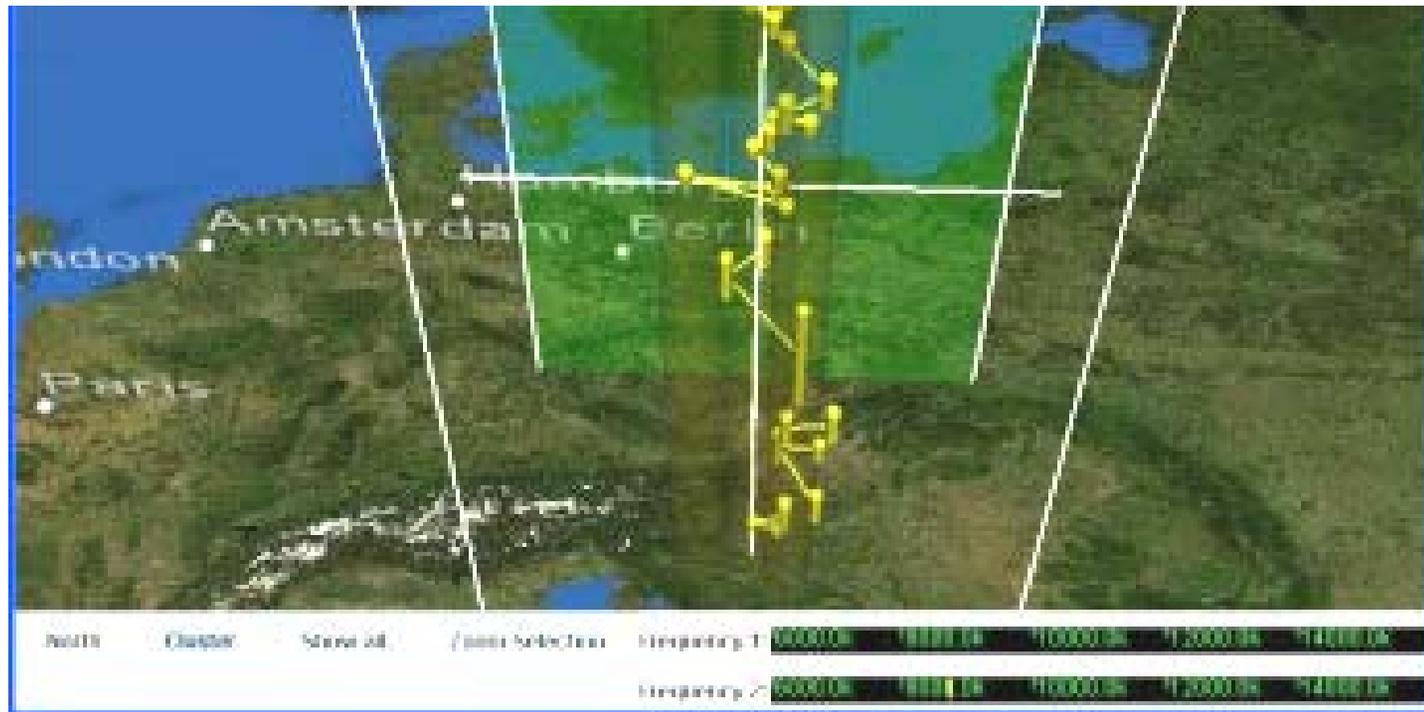
known fixed emitter

We are looking for:

communication partners building a special structure, e.g. A B A C A D A B A C A D ...

# visualisation of a communication structure

## example : free search

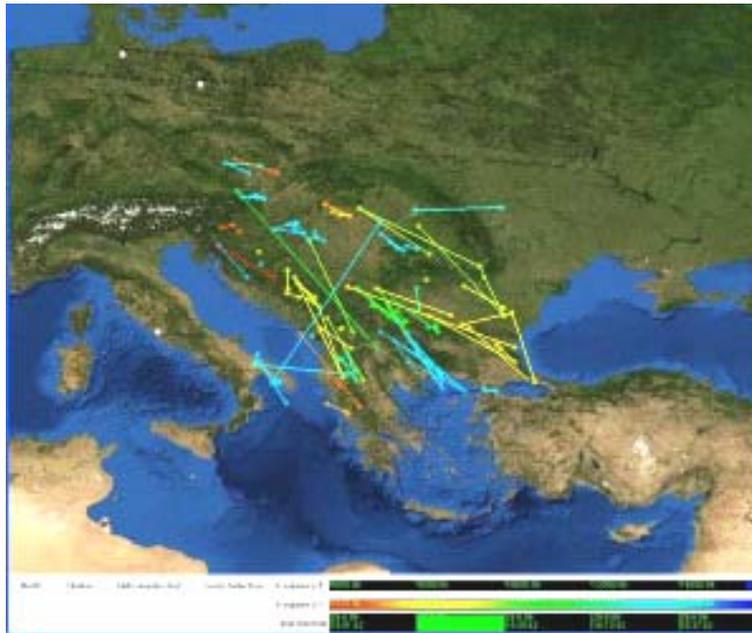


We have: spatio-temporal data selection

We are looking for: two partners communicating in simplex mode

# visualisation of a communication structure

## example: increasing communication activity



6 am ↗



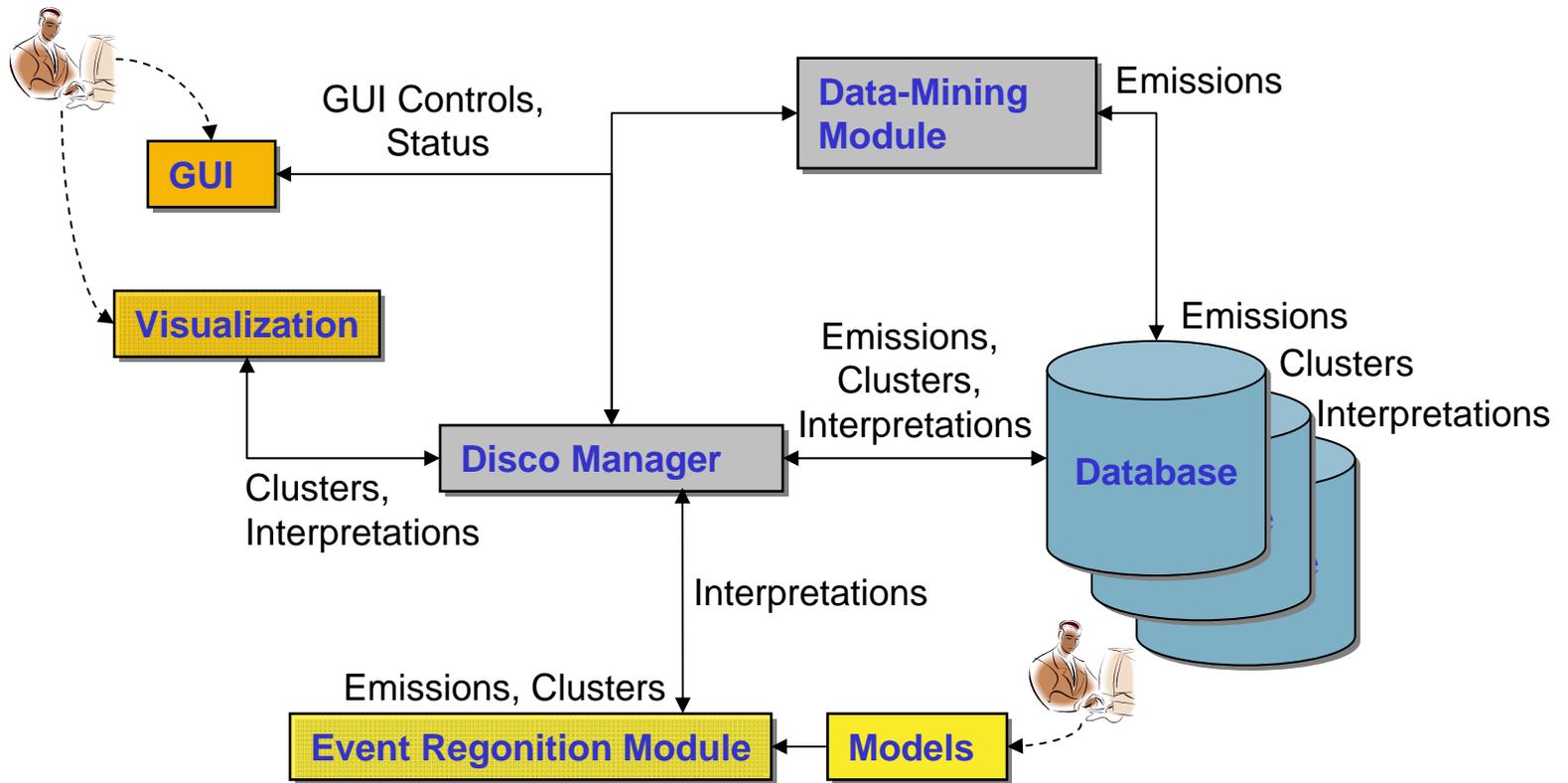
4 pm ↗

At different times the intensity of communication may change heavily.

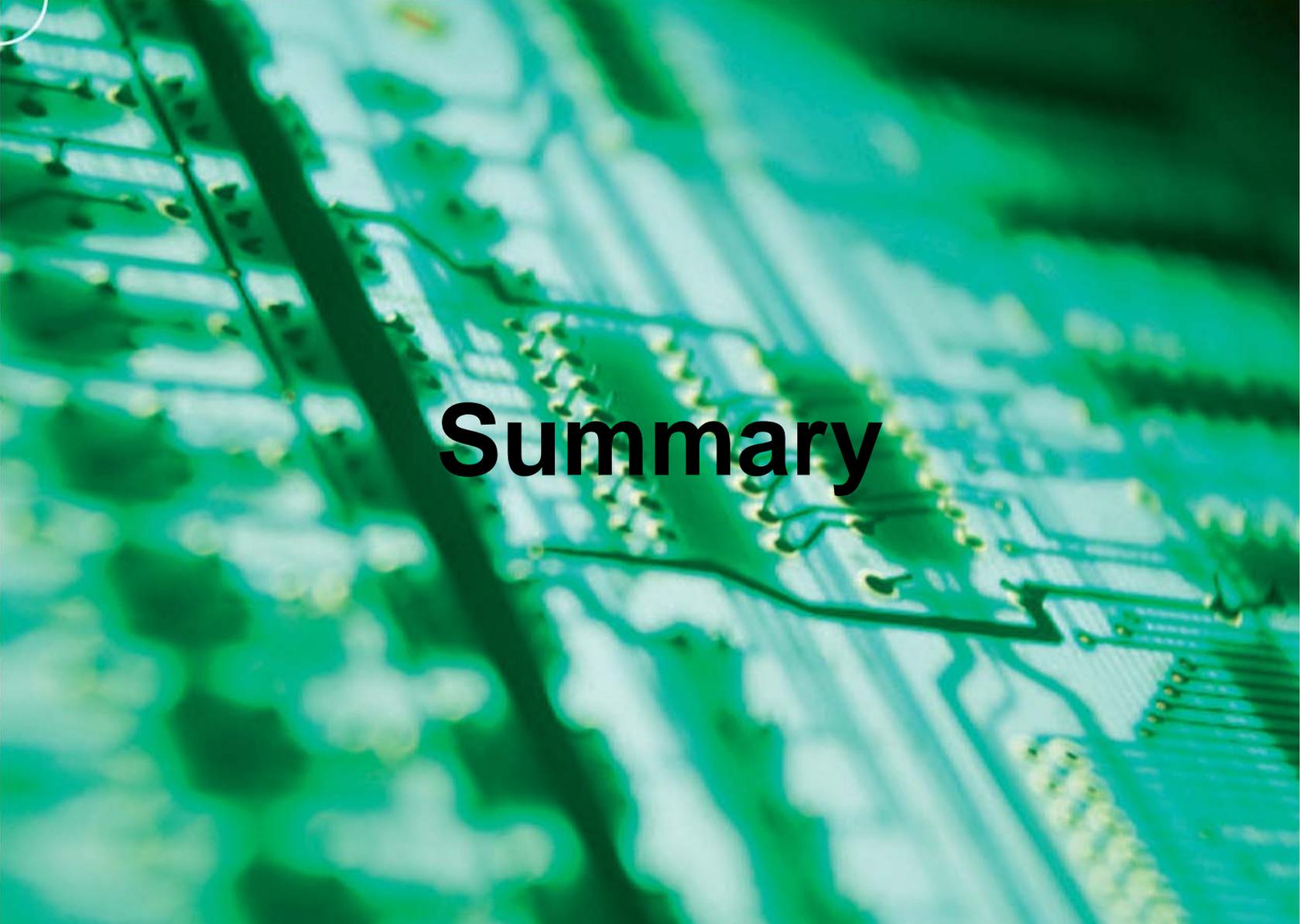
We have: an area of special interest

We are looking for: changes in emission occurrence

# Architecture



existing prototype  
with distributed architecture



# Summary

- Cooperation with computer science research institute competence fields:
  - distributed systems
  - software technique
  - intelligent systems
  - learning with new media
  - logistics simulations
  - usability and software-ergonomics
  - IT-Security
  - visualisation and interactive media
- Incorporation of the newest and best research result
- Know-How and innovation transfer from research into industry

- The COPIN approach makes massive emission data streams manageable and analysable
- Domain specific orientation integrates the expertise of the intelligence analysts
- We aim at meeting the interests of our customers
- COPIN provides an exploration process to successfully add value to information:
  1. data mining
  2. cluster visualisation
  3. model based communication detection
  4. visualisation of communication structures
- The model based concepts represent a flexible independent way to work with your data
- The COPIN concept is applicable also to other data collections.

**Thank you very much  
for your attention!**

**Dr. Vera Kamp  
Dr. Joachim Stamm**

**Tel: +49 40 23734-234  
Fax: +49 40 23734-173  
vera.kamp@plath.de**

**PLATH GmbH**

**Gotenstrasse 18  
20097 Hamburg  
Germany  
www.plath.de**