# EVALUATION OF ARTIFICIAL INTELLIGENCE SYSTEMS

GUILLAUME AVRIN

18/10/2019

# Our positioning in AI

CRÉER LA CONFIANCE

# EVALUATION OF ARTIFICIAL INTELLIGENCE SYSTEMS

**LNE, state-owned trusted third party for the evaluation of AI and robots**

**As a state-owned laboratory:**

It is <u>independent of any private interest</u>
**(reinforced notion of trusted third party)**

**The sincerity of its evaluations is guaranteed**

**More than 10 years of experience on AI evaluation and more than 900 systems evaluated** by a permanent team of doctors and engineers specialized in evaluation.

CRÉER LA CONFIANCE **LNE**

# The evaluation, step by step

CRÉER LA CONFIANCE

# AI EVALUATION PROCESS

# AI EVALUATION PROCESS

**Gray-box evaluation**



Environment

SYSTEM

SENSORS

INTERPRETATION

DECISION

EFFECTORS

Action

Detection evaluation

Interpretation evaluation

Decision making evaluation

Action evaluation

**Black-box evaluation**

Environment

SYSTEM

Action

Overall system evaluation

# ILLUSTRATION OF THE STEPS OF AN EVALUATION

Task definition → Provision of testing datasets and environments → Retrieval of system outputs → Comparison of system outputs and references → Scoring and error analysis

# EVALUATION : AN EXPERTISE IN ITS OWN RIGHT

**Evaluation plan**

**Evaluation references**

| 1. Testing scenarios | 2. Protocols, metrics | 3. Testing environments | 4. Data | 5. References (ground truth) |
|---|---|---|---|---|
| Identification of technoscientific barriers to be removed | Identification of influencing factors | Development of adapted testing environments | Data selection: relevance, representativeness, quality | Development of annotation systems |
| Definition of participation terms and conditions | Definition of evaluation criteria and metrics | Control and measure of influencing factors | Development of tools for data management and sharing (server) | Data annotation or supervision of data annotation |
| Definition of the evaluation tasks | Interpretation of results | Ensure reproducibility of experiments | Development of tools for data collection | Qualification of annotations and annotators |

CRÉER LA CONFIANCE LNE

# LNE EVALUATION TOOLS

| 1. Testing scenarios | 2. Protocols, metrics | 3. Testing environments | 4. Data | 5. References (ground truth) |
|---|---|---|---|---|



Open-source Matics software suite to explore annotated data and evaluation results:

- Translation
- Diarization
- Transcription
- Speaker verification

And soon:

- OCR
- Image recognition





Evaluation of robots:

- laboratory testing (in LNE climatic chambers)
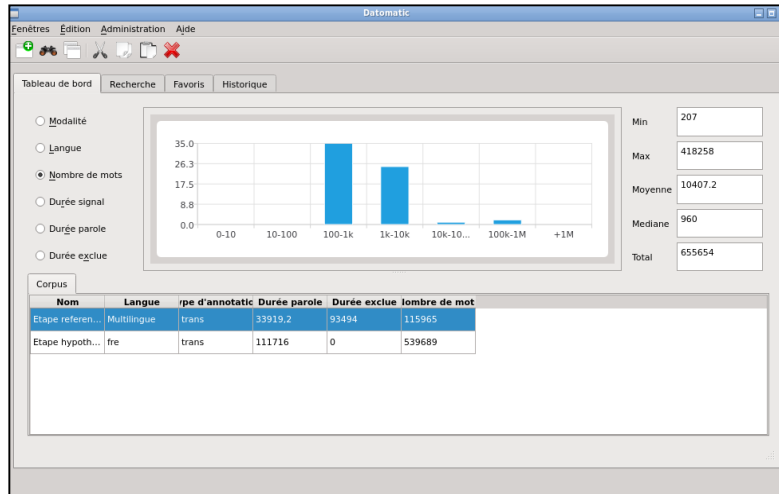- virtual testing (simulation-based)



DIANNE software:

- annotation and automatic pre-annotation of crops and weed
- will be extended to other recognition tasks

CRÉER LA CONFIANCE LNE

# OUR TOOLS

## Matics software suite – Data visualisation and evaluation

Datomatic – Dataset preparation and visualisation
Evalomatic – Evaluation and visualisation



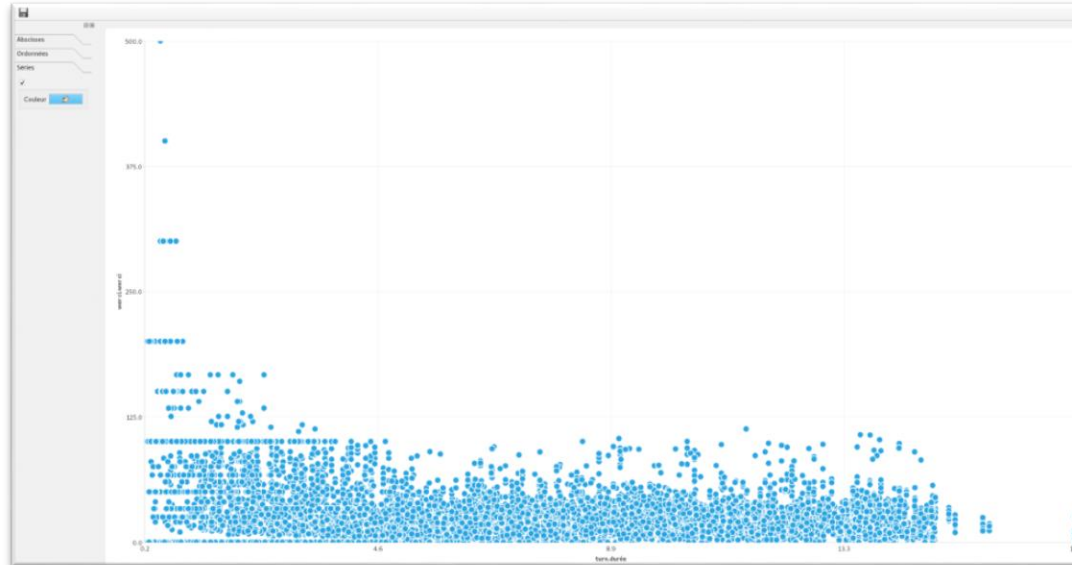***Datomatic***
*Data visualisation*
*(transcription task)*



***Evalomatic***
*Evaluation scores*
*(transcription task)*

# OUR TOOLS

## Matics software suite



***Evalomatic***
*Graphical visualisation*

# OUR TOOLS

## DIANNE : Edge detection, identification and annotation for evaluation

# CHALLENGE ORGANISATION

**Prepa.** — Definition of the evaluation plans and test data or facilities

**Dry-run** — Evaluation protocol validation

**Eval. 1** — First appraisal

**Eval. 2** — Measure of improvements

### Evolution of error rates – person recognition
(REPERE campaign)



Legend:
- Participant 1 supervised
- Participant 2 supervised
- Participant 3 supervised
- Participant 1 unsupervised
- Participant 2 unsupervised
- Participant 3 unsupervised

Y-axis: EGER (%); X-axis: 2012, 2013

### Evolution of error rates – optical character recognition
(MAURDOR campaign)



Legend:
- Participant 1
- Participant 3
- Participant 4

Y-axis: CER (%); X-axis: 2013, 2014

CRÉER LA CONFIANCE — LNE

For which application areas?

# EXPERTISE IN EVALUATION OF INFORMATION PROCESSING SYSTEMS

## SPEECH

Transcription, keyword spotting, speaker comparison, named entities recognition, speaker tracking, translation, etc.

## TEXT

Topic detection, named entities recognition, information retrieval, translation, etc.

ترحيب ، يسعدنا أن نرحب بكم

Welcome, we are delighted to have you here

- Challenges (Quaero, Repere, etc.)
- Benchmarking (INC)
- Qualification (Allies)
- Certification (Voxcrim)

## IMAGE

Head tracking, optical character recognition, etc.

## MULTIMEDIA

Person tracking, document classification, etc.

# EVALUATION OF ROBOTS

- **Smart mobility**     Simulation for autonomous vehicle safety

- **Agri-food**     Risk analysis, scientific monitoring and community structuring, organization of a challenge in agricultural robotics

- **Service**     Development of evaluation tools

- **Public-Private partnership**     Study of the influence of climatic conditions on the performance of AI systems, assessment of AI and cybersecurity of smart medical devices

Simulation of the autonomous vehicle

HRP2 robot (Franco-Japanese humanoid robot) evaluated in climatic chambers at LNE

CRÉER LA CONFIANCE   LNE

# Our orientations

## OUR ORIENTATIONS

**Metrology:** develop standards and protocols for the evaluation of AI

**Evaluation:** set up an AI assessment and testing centres

**Certification:** promote the certification of AI

# METROLOGY OF AI

**Definition of standards:** reference testing datasets and environments, metrics, etc.

**Definition of evaluation protocols:** testing scenarii, evaluation tasks, methods for calculating the measurement uncertainty, etc.

**For performance evaluation**

- Accuracy, precision, trueness, fidelity, error rate, sensitivity, specificity, etc.
- Robustness, resilience and operating range
- Datasets qualification (representativeness)
- Other performance requirements (speed, efficiency, ergonomics, etc.)

**To promote acceptability**

- Regulation (transparency, non-discrimination)
- Explainability, intelligibility, predictability, readable behaviour
- Security (controllable, auditable)

CRÉER LA CONFIANCE LNE

# EXPLAINABILITY

**A tool to facilitate verifications and make them more reliable**

- Solving the "black box" problem?
- To estimate the operating domain, better identify rare (but critical) phenomena, etc.

**Towards the evaluation of explainability**

- Measuring performance
  - Characterise explainability (according to context, requirements, user profile, etc.)
  - Define objective metrics
- Development of standards
  - Type of information to be extracted, reference values, etc.

CRÉER LA CONFIANCE LNE

# CONCLUSION

## What LNE offers:

- A unique know-how in the organization of evaluation campaigns for AI systems (design of the evaluation plan, organization of evaluation meetings, management of the associated events)
  - to set up a rigorous metrological approach (repeatable performance measurements, reproducible experiments, qualified test databases, identified and controlled influence factors, limited biases)
  - to maximize the impact of evaluations
- Evaluation tools
  - Suite Matics software suite
  - annotation tools
  - real or simulated test environments, etc.
- A status:
  - trusted third party (LNE does not develop AI systems)
  - independent evaluator (LNE is public, it is independent of any private interest)

## LNE is interested in:

- collaborating with other **NMI** to bring metrology expertise to the field of AI evaluation.
- participating in projects aimed at demonstrating the performance and functionality of AI technologies
- setting up challenges, evaluation campaigns, competitions, especially in AI and robotics

# COLLABORATIVE TOPICS WITH LNE

**Performance evaluation:** accuracy, precision, trueness, robustness, resilience
- at the level of the overall system (autonomous cars, surgical robots, etc.),
- at the level of the detection modules (obstacle detection, face recognition, etc.),
- at the level of decision-making modules (hazard management, etc.),
- At the level of action modules (autonomous navigation, etc.).

**Explainability evaluation:** how to relate the decision taken to the known data and characteristics of the situation?

**Human-machine interaction evaluation:** how to measure the quality of an interaction (during a close cooperation between an intelligent personnel assistant and a pilot, for example).

# Thank you for your attention

CRÉER LA CONFIANCE LNE