# LCT ANNUAL MEETING 2018

## CONFERENCE PROGRAM

## NANCY, FRANCE
## JUNE 27-29

Co-funded by the
Erasmus+ Programme
of the European Union

UNIVERSITÉ
DE LORRAINE

Loria
Laboratoire lorrain de recherche
en informatique et ses applications

# Contents

# Preface

Dear LCT students, partners and participants,

It is our great pleasure to welcome you all to the LCT Annual Meeting. This year we will celebrate the 12th anniversary of the LCT consortium and  the return of the LCT Annual meeting to Nancy.

The Annual Meeting  gathers all students and scholars participating in the LCT Masters' program, and  provides the opportunity for students to share their experiences, to present their latest research results and to showcase their culture.  It will also offer the opportunity to discuss several organisational aspects of LCT. The program includes the participation of  3 invited speakers that will give keynote talks of different nature and flavour, notably,

- Barbara Plank (University of Copenhagen) *Learning across tasks and languages*
- Gerard Casanova (University of Lorraine) *Soft skills for employability: main outputs of a European project*
- Aurélie Névéol (LIMSI, CNRS) *Natural Language Processing for epidemiology and public health*

In this booklet you will find further information on the invited talks and speakers, the abstracts of the contributed posters, the detailed program of the Annual Meeting, as well as some useful information that we hope that will improve your stay.

We would like to thank the President of the University of Lorraine (Pierre Mutzenhardt), the director of LORIA (Jean-Yves Marion), and the keynote speakers for kindly accepting our invitation to participate in this year's Annual Meeting. I would also like to personally express my deepest gratitude to all that contributed to the organization of our event, namely,  the student organization committee (Gosse Minnema, Maria Andrea Cruz Blandon, Aria Nourbakhsh), Olivia Brenner (LORIA) and Pascale Lacroix-Malgras (IDMC, University of Lorraine). Finally, I would like to thank LORIA, IDMC (former UFR Math-Info) and of course the consortium Erasmus Mundus LCT, for the financial and practical support.

We are looking forward to meeting you all in Nancy! We hope that you will enjoy your stay in Nancy and  that this year's Annual Meeting will leave in each one of you excellent memories!

Miguel Couceiro

Local LCT Coordinator, University of Lorraine

# Invited speakers

**Gérard Casanova (eLene4work)**

*Soft skills for employability*

*Gérard is a lecturer at the Université de Lorraine and currently works in the digital utilisation department* (direction des usages du numérique). *He has been actively involved in digital learning for more than 15 years where the focus of his work has been on the creation and tutoring of online training courses and the support and training of those responsible for the design and production of learning resources as well as tutoring at a distance. Gérard's work has included pedagogical responsibility for a full online masters course. He worked in several European projects such as eLene4work.*

The presentation is based on an European project (*elene4work*). The first part of the presentation is about the goals of the eLene4work project as well as the partners of the project. In a second part will be presented a comparative analysis on the state of the art of soft skills and digital soft skills including both methodology and main results. Then we will speak about which soft skills students have and which should they have? What is a soft skill? Then I will present the eLene4work soft skills framework and the learning path proposed by the elene4work project: a self-assessment tool, orientation guide and the personal journal. At the end I will present the lessons learnt from the project and the following of the project, at least in France.

**Talk** Day 2, 09:00 | Venue B

## Aurélie Névéol (LIMSI/CNRS)
*Natural Language Processing for epidemiology and public health*

*Aurélie Névéol is a Senior Staff Scientist at the* Centre National pour la Recherche Scientifique *(CNRS, LIMSI). She leads research on clinical natural language processing for languages other than English. Her research work includes using NLP to create representations of clinical information to support information extraction from unstructured clinical narrative text in the electronic health record, which can then be used for high throughput phenotyping. She earned an MSc in Linguistics from Université Paris 7 in 2002 and a PhD in Computer Science from INSA de Rouen in 2005. She has then contributed to research projects at the National Library of Medicine to improve the retrieval and analysis of biomedical text from the litterature. Dr. Névéol has also been the primary organizer of the information extraction task at the CLEF eHealth lab since 2015.*

Much clinical and biomedical knowledge is contained in the text of published articles, Electronic Health Records (EHRs) or online patient forums and is not directly accessible for automatic computation. The goal of biomedical Natural Language Processing (bioNLP) is to extract information from free text narratives and convert it to machine-readable representations that can be integrated in clinical workflows. In recent years, there has been a growing interest in leveraging bioNLP for applications in epidemiology and public health research. In this talk, I will present current methods of bioNLP in English as well as other languages, illustrate the implications of this research work in clinical applications and outline research challenges that still need to be tackled.

**Talk** Day 2, 14:00 | Venue B

**Barbara Plank (IT University of Copenhagen)**
*Learning Across Tasks And Languages*

*Barbara Plank is Associate Professor in Natural Language Processing at IT University of Copenhagen. She has previously held positions as tenured assistant professor at the University of Groningen, as assistant professor and postdoc at the University of Copenhagen and as postdoc at the University of Trento. After finishing her Master's degree in EM LCT master in 2007 (University of Amsterdam and Free University of Bozen-Bolzano), she pursued a PhD at the University of Groningen. Her research interests include learning under sample selection bias (domain adaptation, transfer learning), annotation bias and generally, semi-supervised and weakly-supervised learning (learning under limited supervision) for cross-domain and cross-lingual NLP.*

How can we build Natural Language Processing models for new tasks and new languages?  In this talk I will survey some recent advances to address this challenge, from multi-task learning, cross-lingual transfer to learning multilingual models and learning models under tiny supervision.

**Talk** Day 1, 11:00 | Venue B

# Student posters

## Badr Abdullah (University of Lorraine)
*Leveraging Word Contexts in Wikipedia for Recovering OOV Proper Nouns in Speech Recognition*

Automatic Speech Recognition (ASR) systems are usually trained on static data and a finite word vocabulary. When a spoken utterance contains Out-Of-Vocabulary (OOV) words, ASR systems misrecognize these words as in-vocabulary words with similar acoustic properties, but with entirely different meaning. The majority of OOV words are information-rich proper nouns (e.g., person names, geographic locations, commercial products) that are vital to spoken content understanding. Therefore, failing to recognise OOV words has a significant adverse impact on many downstream applications such as spoken document indexing and translation.

In this thesis, we address this problem by dynamically extending the ASR vocabulary based on the context obtained from ASR initial first-pass hypothesis. In other words, given the in-vocabulary transcription of a spoken utterance, the goal is to retrieve a ranked list of OOV proper nouns that are likely relevant to the context, add words in this list to the ASR lexicon, and perform a second-pass decoding with the ASR system. To this end, we explore different techniques that leverage topical contexts of OOV words in Wikipedia to develop neural models for OOV words retrieval.

## Xiaoyu Bai (University of Groningen)
*Automatic Classification of English Learner Proficiency Using Elicited Versus Spontaneous Data*

Automatic classification of language learner texts to proficiency levels can help learners in terms of self-assessment and help teachers determine their approximate level. For this task, data is typically drawn from foreign language exams, which have been produced in a regulated context in response to specific prompts. While such data are clean, they are rare, often difficult to access and do

not represent language production in a natural setting. Therefore, we perform in this study the automatic classification of non-native English speakers' proficiency levels using spontaneous text production on social media websites, as compared to using data in a language teaching context.

On the one hand, we harvest data from Twitter and Reddit in which non-native English speakers report on their own proficiency levels using phrases such as "I'm at B2 in English" and treat such self-reported levels as a proxy to their true proficiency labels. We then extract the accessible posts by the same authors written in English and obtain a proficiency-labeled dataset of social media English data. On the other hand, we use data drawn from the *EF-Cambridge Open Language Database (EFCamDat)* corpus (Huang et al., 2018; Geertzen et al., 2013), which contains learner writings from the online English course *EF English Live*, grouped into multiple proficiency levels.

We carry out classification and evaluate system performance in both domains with the aim of examining overlaps, contrasts and possible interactions between the predictive features in the two domains. We also discuss the suitability of using social media data as spontaneous language production for this task and its potential contribution to improving proficiency level classification in general.

**Angelo Basile, Kenny Lino & Urana Urosevic (University of Malta)**
*Textual classification of gay men: a linguistic and ethical experiment*

"Finally, the ethics of classifying. Which tools should we not build?" (Bender, 2018).

In this work we want to classify gay men through their texts. We leverage publicly available data from online forums such as Reddit and Twitter: we use distant supervision to collect a large amount of text and use a variation of a state-of-the-art gender prediction system to classify whether a certain author is gay or not. Also, to minimize the topic bias, we make sure to filter the data throughout the process.

Previous work in linguistics has shown that text can be used to predict social variables such as gender (Rickford et. al., 2013), age (Barke, 2000), location

(Wieling et al., 2011) and even Myer-Briggs personality type (Verhoeven et al., 2016). However, work investigating the interaction of sexual orientation and writing style seems to be far and few. This can be partially attributed to the lack of easily available labeled data - in fact, an examination of common author profiling corpora show that only the CSI Corpus (Verhoeven and Daelemans, 2014) contains such information. Following the idea that language variation is influenced by social settings and variables (Eckert, 2008), we define gay men here to be those who index themselves as such in their writing. That is to say, in order to be considered gay in this investigation, we require that the authors both identify as gay themselves and present that identity in their work.

Of course, such classification is hazardous. Through this work we present a computational sociolinguistic experiment, raise awareness on how publicly available texts/data can be tapped for all kinds of NLP works and purposes, and contribute to the current hot debate on data processing ethical issues.

## Somaye Jafari Tazehjani (University of Malta)
*From Entailment to Generation*

Natural Language Understanding (NLU) is a challenging task with regard to the high level of complexity of natural languages and it has been of long-standing interest in different scopes, such as, computational semantics and Natural Language Processing (NLP). In earlier works, the ability to reason with natural language was gained through different logical and statistical approaches, whilst, recently, there has been a growing interest in the intersection of these two fields, NLP and, Computer Vision (CV). In this approach, researchers are interested in multimodal systems and feeding the models with visually-augmented data rather than only unimodal systems and inputs.

Recognizing Textual Entailment (RTE) which is a necessary step towards true NLU is the prerequisite to many NLP tasks, such as summarization, question answering, and Machine Translation (MT) systems. In a recent work, which RTE was framed as a classification task, taking the approach of adding visual information to the input showed improvement in results compared to taking only unimodal input. In the current work, on other hand, textual entailment problem is dealt with as a generation task, which is a useful approach, considering the fact that in many cases, only premises are available, not premise, hypothesis pairs.

To deal with the entailment generation task, we implement sequence-to-sequence architecture for neural machine translation (NMT) calling the source and target language 'encoder' and 'decoder' respectively. Moreover, we investigate whether the performance of the model improves if, in addition to premise and hypothesis, there is also an image which corresponds to the relevant premise. For the purpose of training,We use an entailment-subset of multimodal version of SNLI, in which the premises are mapped onto the original image, from Flickr30k dataset. Finally, to evaluate the generated output, the same metrics as MT evaluation metrics such as BLEU and METEOR are utilized.

## Dan Kondratyuk (Charles University in Prague)
*Beating State-Of-The-Art in Czech Part-Of-Speech Tagging and Lemmatization*

We present LemmaTag, a featureless recurrent neural network architecture that jointly generates part-of-speech tags and lemmatizes sentences of languages with complex morphology, using bidirectional RNNs with character-level and word-level embeddings. We demonstrate that both tasks benefit from sharing the encoding part of the network and from using the tagger output as an input to the lemmatizer. We evaluate our model across several morphologically-rich languages, surpassing state-of-the-art accuracy in both part-of-speech tagging and lemmatization in Czech, German, and Arabic.

## Ajinkya Kulkarni (University of Lorraine)
*Expressive speech synthesis based on transfer of prosodic information*

Current interactions with Human-machine interface is still constrained with more or less natural voice without consideration of the relationship between textual semantics and emotions in speech. The current speech style is typically a "reading style", which results from the style of the speech data used to develop TTS systems. To create expressive speech corpus is a complex process in terms of accuracy of expressions, time involved in development and resources invested. In order to reuse the existing data on different users, we need expressivity adaptation techniques.

The focus of work lies on acoustic and duration feature space transformation from neutral to expressive for TTS system. This mechanism will allow avoiding recording of additional expressive speech corpus, thus providing an ability to generate expressive speech synthesis system. We propose conditional variational autoencoders to learn the latent representation of speech conditioned on text. Furthermore, we demonstrate the ability of model to transfer the expressivity to neutral speech content, without need of parallel speech data.

**Aria Nourbaksh, María Andrea Cruz Blandón & Gosse Minnema (University of Lorraine)**
*What's The Answer: Dialogue Annotation*

Our project focused on identifying and classifying question-answer pairs in spoken language. We had three main goals: (1) designing an annotation schema for classifying questions and answers, (2) writing an annotation guide and manually annotate dialogues in different languages, and (3) exploring machine learning approaches to automate our work. Our work took place in the context of the SLAM (Schizophrénie et Langage: Analyse et Modélisation) and aims to contribute to the development of computational models of both impaired and unimpaired discourse.

Our annotation schema comprises five types of questions classified based on their form (syntax) and function (semantics and pragmatics) and seven types of answers based on the possible answers a particular type of question can have. To test this schema, we wrote an annotation guide that describes the tagging process.
We then used this annotation guide to tag dialogues in three different languages: English, Spanish, and Dutch. A part of the English corpus was designated as the 'golden standard' and used to compute agreement scores. Agreement was moderate for question types ( $\kappa = 0.63$ ), but less so for answer types ( $\kappa = 0.49$ ). Since our annotation task is quite complicated and requires making often subtle distinctions, this was only to be expected.

We experimented with two approaches for automatic question type classification: a classical statistical approach (decision trees) and several neural network approaches (bag-of-words classifier, RNN classifier). The decision tree produced the best results (accuracy 73% , F1 = 0.58 on unseen data); the neural networks generally performed well on seen data but were unable to generalize. All models suffered from a lack of sufficient training data (only ~200 data points as input);

until more human annotations are produced, it is too early to draw any definitive conclusions.

## Adedayo Oluokun (Charles University in Prague)
*Creation of a Dependency Treebank for Yoruba using Parallel Data*

The goal of this research is to create a dependency treebank for Yoruba, a language with very little pre-existing machine-readable resources. The treebank will follow the Universal Dependencies annotation standard. Known techniques for porting resources from resource-rich languages will be tested, in particular projection of annotation across parallel bilingual data. Manual annotation is not the main focus of this research; nevertheless, a small portion of the data will be verified manually in order to evaluate the annotation quality.

## Thuong-Hai Pham (Charles University in Prague)
*Exploiting Sentence Structure in Neural Machine Translation*

Neural machine translation (NMT) has become the new state of the art. Similarly to the previous best-performing approach of phrase-based MT, NMT is linguistically uninformed. Our goal is to focus on the state-of-the-art Transformer model and attempt to promote the knowledge of source-side syntax either by enriching the encoder with linguistic information about the sentence structure or by multi-task learning. We propose a novel idea to interpret Transformer self-attention as a dependency parse, hence, the model can learn to jointly translate and parse the source simultaneously. Our results show that a variety of enriching methods improved translation quality, yet not significantly. However, the treatment of self-attention as dependencies shows improvement in translation and good performance in parsing. With minimal modification of the Transformer, our joint model suggests that the Transformer model can very easily grasp this structure.

## Josine Rawee (University of Trento)
*Visual features in distributional semantics*

Distributional semantic models have been found to be able to represent the meaning of words based on co-occurrences in big text corpora. These models can

approximate semantic similarity as tagged by human annotators. However, being able to extract visual properties, such as the prototypical color of a concept, from distributional models seems to be very hard. This is usually attempted by comparing vectors of concepts with vectors of colors. We build instead a semantic space that is dependent on a count model where target words are only defined by their co-occurrences with color terms. We will explore for different categories the distribution of colors found, to understand how these concepts are described in terms of their color. Our hypothesis is that distributional spaces, as representations of conceptual knowledge, can deviate quite sharply from what the world is like (or even what speakers think the world is like). For example, many plants are green, but is green a salient property of plants in a semantic space? Our investigation attempts to clarify when language use matches speakers' perception of the world, and when it doesn't.

**Anastasia Serebryannikova (University of Groningen)**
*Predicting Stock Market Trends from News Articles*

The goal of predicting stock market trends has emerged from the desire to make profitable investing decisions. Random Walk Theory and Efficient Market Hypothesis suggest that it is impossible to outwit the market and the stock prices are changing at random. However, recent advances in machine learning and the growing availability of wide-scale data have made it possible to apply state-of-the-art algorithms to this problem.

We will make an attempt to analyze financial news articles in order to predict the changes in the stock prices. Even though a lot of work has been conducted in this area, there are still many weak points. Our research will address some of them.

The first issue concerns the evaluation standards. The existing evaluation methods focus on different subtasks; various time frames are chosen for the prediction and different observation windows are used. Therefore, even if the textual data as well as the stock market information are available, it is not self-evident how to align them in order to get meaningful results. We will explore the possible data alignment options and find the ones yielding the best performance.

Second, most researchers used different datasets for solving this task and therefore it is impossible to compare the implemented approaches directly. If

feasible, we will try to compare some of the existing approaches on the same dataset and explore the impact of various textual, meta-textual and non-textual features that were mentioned in the previous works and that we come up with ourselves.

Last but not the least, the performance of the best classifiers never exceeds 65% for the binary classification task focusing on the directionality of the change, which obviously shows that there is still some room for improvement. We will try to achieve better results by combining the findings from previous research with our own conclusions.

## Svetlana Tchistiakova (University of Trento)
*Acoustic Models for Second Language Learners*

Automatic speech recognition (ASR) systems have been slow to become adopted in the field of language teaching, in part because ASR systems still face difficulties when encountering speech from non-native speakers, child speakers, and noisy environments-- the prevailing conditions in the classroom. In addition, well-annotated, spontaneous speech data from native language (L1) speakers of a particular second language (L2) are not always plentiful, particularly from classroom environments. Furthermore, when available, this type of speech data exhibits difficult-to-model L2 phenomena including phonological crossover (i.e. accent), lexical crossover, pronunciation errors, hesitations, disfluencies, and background noise. We show that during training of an ASR system, by including relatively little in-domain multilingual L2 speech data, adding in spontaneous L2 speech data, and including extra examples of background noise, we are able to improve word error rate (WER) performance by approximately 40% absolute.

## Ludmila Tydlitátová (University of Groningen)
*Automatic Extraction of Adverse Events from Literature*

Drug safety (pharmacovigilance) is a broad term that describes the collection, analysis, monitoring and prevention of adverse events in drugs and therapies. Adverse events are monitored in several ways, for example through reports written by medical experts, through social media, or by studying academic journals that

report on case studies or experiments. Monitoring journals for mentions of specific drugs and possible adverse events is a highly regulated process that is typically conducted manually by trained professionals. We will explore and evaluate statistical and machine learning methods that can help automatically scan academic literature for essential knowledge about potential adverse events. Previously some analysis has been carried out on English data. We will be working with articles written in Czech and Slovak.

## Claudia Zaghi (University of Groningen)
*Hate speech detection in social media*

Hate speech is "the use of aggressive, hatred or offensive language, targeting a specific group of people sharing a common trait: their gender, ethnic group, race, religion, sexual orientation, or disability".

The phenomenon is widely spread online. To monitor the problem, social networks and websites have introduced a progressively stricter code of conduct and regularly remove hateful content flagged by users. However, the volume of data in social media makes it challenging to supervise the published content across platforms and languages.

Automatic detection of hate speech has therefore been an active task in Natural Language Processing in the last years, addressed via resource creation, and detection systems, exploiting different combinations of surface, lexical and morphosyntactic features. Currently, the majority of systems are supervised, thus requiring the manual annotation of training data, and are tailored for one language only, with very limited portability.

The research we present addresses the issue of hate speech in Italian and German. The aim of the study is to automatically detect hateful content, based on data scraped from different social media (Facebook, Twitter, and YouTube) and the resources made available by the 2018 Germeval and Evalita shared tasks on hate speech. Dealing with two languages and different platforms at once we hope will highlight potentially common features that could be exploited to enhance portability and/or ease the development of new resources and systems. Specifically, we will describe the existing resources and those we have built anew, and report results on three aspects: (i) the validity of using distant supervision to

acquire silver training data, thus reducing the need for manual annotation; (ii) the contribution of different methods of acquiring lexical resources that can be used in classification; (iii) the study of the effectiveness of the developed machine learning models across platforms and languages.

## Other student posters

*Mostafa Abdou (University of Groningen)*
*Reyhaneh Hashempour (University of Malta)*
*Artur Kulmizev (University of the Basque Country)*
*Guido Linders (University of Trento)*
*Daniel Low (University of Groningen)*
*Vinit Ravishankar (Charles University in Prague)*

# Venues

**Loria (615 Rue du Jardin Botanique, 54506 Vandoeuvre-lès-Nancy)**
*Access: tram #1, stop "Callot" (from city center: direction "Vandoeuvre CHU Brabois")*
*+ 5 mins walking (walk past U Express, turn right towards Lycée Jacques Callot, then left*
*until you reach the Campus Sciences. Go through the barrier; Loria is across the parking*
*lot, to the left of the big orange University Library building). N.B.: entering the lab can*
*take some time due to anti-terrorism measures.*

Venue A: room A006 (Salle Jacques Louis Lions)
Venue B: room A008 (Salle Jean Legras)
Venue C: room B013 (Salle Bob)

# Student accomodation

**CROUS Résidence Boudonville (61 Rue de Boudonville, 54052 Nancy)**

*Access: bus #2, stop "Aimé Morot" + 3 mins walking*

*To city centre: bus #2, direction "Laneuville Centre" (get off after "Nancy Gare")*

*To Loria: tram #1 from "Nancy Gare", direction "Vandoeuvre CHU Brabois"*

# Nancy City Centre

*Access: tram #1, stop "Point Central" or "Cathédrale" / bus #2, stop "Place Stanislas"*

**What to do in Nancy?**
Place Stanislas is the best-known place of Nancy, you should definitely not miss this beautiful square! There, you will find the Museum of Fine Arts, The Opera House, the park La Pépinière and plenty of options for enjoying French cuisine and wines. If you are an ice-cream lover, do not forget to go to "Amorino" (just follow Stanislas' pointing finger!).

Palais du Gouvernement

Kiosque Mozart

Horloge Florale

Auditorium Gaston Stoltz

Allée Léon Tonnelier

Hôtel Héré

P

Place de la Carrière

Place de la Carrière

Rue des Écuries

Mini-Golf

Bâtiment H

Cité Administrative

P

Buste de Sellier

Rue du Maure Qui Trompe

Grande Rue

Le Tém'

Terrasse de la Pépinière

Bâtiment Y

Cour d'Appel Tribunal de Grande Instance

Le Romana

Rue Callot

Le Crêp'show

Statue Jacques Callot

La Cerise sur le Gâteau

L'Arrosoir

Place Nelson Mandela

Opéra Café

P Parking Vaudémont

Rue Sainte-C

Rue Guibal

Rue Lyautey

Rue Sainte-Catherine

La Maison dans le Parc

L'Alsace à Table

Le Potager

Papa Joe

Suzette

Fontaine Amphitrite

Glaces Lorraines

Jean Lamour

Bastion d'Haussonville

Jardin du musée des Beaux-Arts

Place d'Alliance

Fontaine de Guibal

Musée des Beaux-Arts

Statue Stanislas Leszczyński

Grand Hôtel de la Reine

Préfecture de Meurthe-et-Moselle

Rue Pierre Fourier

Con de la Re de L

Qui l'eût cru

Grand Café Foy

Les Portes d'Or

Rue Gambetta

Sushi Shop

Hôtel de ville

Rue Pierre Fourier

Rue Saint-Julien

Lycée Jeanne d'Arc

Rue des Dominicains

Rue Claude Charles

Rue Maurice Barrès

P Parking Place Stanislas

Banque Kolb

Rue Dom

Cathédrale

Parvis

21

# Schedule Day 0

| TBA | Partner Welcome & Dinner    Student Dinner & Bowling |
|-----|-----|

# Schedule Day 1

| 8:30–9:00 | **Venue B**<br>Welcome |
|-----------|-----------|
| 9:00–9:15 | **Venue B**<br>Introduction by Ivana Kruiff-Korbayova |
| 9:15–10:30 | **Venue A**<br>Consortium Meeting    **Venue B**<br>Student Meeting |
| 10:30–11:00 | **Venue C**<br>Coffee Break |
| 11:00–12:00 | **Venue B**<br>Barbara Plank: *Learning Across Tasks And Languages* |
| 12:00–13:30 | **Venue C**<br>Lunch |
| 13:30–14:15 | **Venue B**<br>Official welcome (with the president of the University of Lorraine, the director of LORIA, and the director of the UFR Math-Info) |
| 14:15–16:00 | **Venue B**<br>Students & Coordinators (+ Coffee Break **15:00-15:15**) |
| 16:00–17:30 | **Venue B**<br>Graduation Ceremony (with the head of the NLP Master's Program in Nancy and the LCT Student Representative) |
| 17:30–18:30 | **Venue C**<br>Appetizers & Drinks |
| 20:00 | Dinner |

# Schedule Day 2

| | |
|---|---|
| **9:00-10:00** | **Venue B** <br> Gérard Casanova: *Soft skills for employability* |
| **10:00-10:30** | **Venue B** <br> Discussions |
| **10:30-12:30** | **Venue C** <br> Poster Session |
| **12:30-14:00** | **Venue C** <br> Lunch |
| **14:00-15:00** | **Venue B** <br> Aurélie Névéol |
| **15:00-16:30** | **Venue C** <br> Closing Meeting |
| **16:30-18:30** | Nancy Walking Tour |