



## The 65th ASH Annual Meeting Abstracts

## ORAL ABSTRACTS

## 803. EMERGING TOOLS, TECHNIQUES AND ARTIFICIAL INTELLIGENCE IN HEMATOLOGY

**Clinical Text Reports to Stratify Patients Affected with Myeloid Neoplasms Using Natural Language Processing**

Gianluca Asti, MSc<sup>1</sup>, Elisabetta Sauta, PhD<sup>1</sup>, Nico Curti, PhD<sup>2,3</sup>, Gianluca Carlini, PhD<sup>2,3</sup>, Lorenzo Dall'Olio, PhD<sup>3,2</sup>, Luca Lanino, MD<sup>4</sup>, Giulia Maggioni, MD<sup>1</sup>, Alessia Campagna, MD<sup>1</sup>, Marta Ubezio, MD<sup>1</sup>, Antonio Russo, MD<sup>1</sup>, Gabriele Todisco, MD<sup>1</sup>, Cristina Astrid Tentori, MD<sup>1</sup>, Pierandrea Morandini, MEng<sup>1</sup>, Marilena Bicchieri, PhD<sup>1</sup>, Maria Chiara Grondelli, BSc<sup>1</sup>, Matteo Zampini, PhD<sup>1</sup>, Erica Travaglini, BS<sup>1</sup>, Victor Savevski, MEng<sup>1</sup>, Nicolas Riccardo Derus, PhD<sup>5</sup>, Daniele Dall'Olio, PhD<sup>5</sup>, Claudia Sala, PhD<sup>5</sup>, Lin-Pierre Zhao, MD<sup>6</sup>, Armando Santoro, MD<sup>1</sup>, Shahram Kordasti, MDPhD<sup>7</sup>, Valeria Santini, MD<sup>8</sup>, Anne Sophie Kubasch, MD<sup>9</sup>, Uwe Platzbecker, MD<sup>10</sup>, Maria Diez-Campelo, MD PhD<sup>11</sup>, Pierre Fenaux, MD PhD<sup>6</sup>, Amer M. Zeidan, MBBS, MHS<sup>12</sup>, Torsten Haferlach, MD PhD<sup>13</sup>, Gastone Castellani, PhD<sup>5</sup>, Matteo Giovanni Della Porta, MD<sup>1</sup>, Saverio D'Amico, MSc<sup>14</sup>

<sup>1</sup>Humanitas Clinical and Research Center, IRCCS, Rozzano, Italy

<sup>2</sup>Department of Physics and Astronomy, University of Bologna, Bologna, Italy

<sup>3</sup>Data Science and Bioinformatics Laboratory, IRCCS Institute of Neurological Sciences of Bologna, Bologna, Italy

<sup>4</sup>Humanitas Clinical and Research Center, IRCCS, Rozzano, Italy

<sup>5</sup>University of Bologna, Bologna, Italy

<sup>6</sup>Saint Louis hospital AHP, Paris, France

<sup>7</sup>King's College London, London, United Kingdom

<sup>8</sup>MDS Unit, DMSC, AOU Careggi, University of Florence, Firenze, Italy

<sup>9</sup>Department of Hematology, Cellular Therapy, Hemostaseology and Infectious Diseases, University Medical Center Leipzig, Leipzig, Germany

<sup>10</sup>Department of Hematology, Cellular Therapy, Hemostaseology and Infectious Diseases, University Leipzig Medical Center, Leipzig, Germany

<sup>11</sup>Department of Hematology, Salamanca-IBSAL University Hospital, Salamanca, Spain

<sup>12</sup>Section of Hematology, Department of Internal Medicine, Yale School of Medicine and Yale Cancer Center, New Haven, CT

<sup>13</sup>MLL Munich Leukemia Laboratory, Munich, Germany

<sup>14</sup>Humanitas Clinical and Research Center, IRCCS, Rozzano (Milan), Italy

**Background:** The availability of multimodal patient data, such as demographics, clinical, imaging, treatment, quality of life, outcomes and wearables data, as well as genome sequencing, have paved the way for the development of multimodal clinical solutions that introduce personalized or precision medicine. The clinical report is an information layer that contains relevant information about the disease in addition to the patient's point of view. Natural language processing (NLP) is a branch of artificial intelligence (AI) and its pre-trained language models are the key technology for extracting value from this data layer.

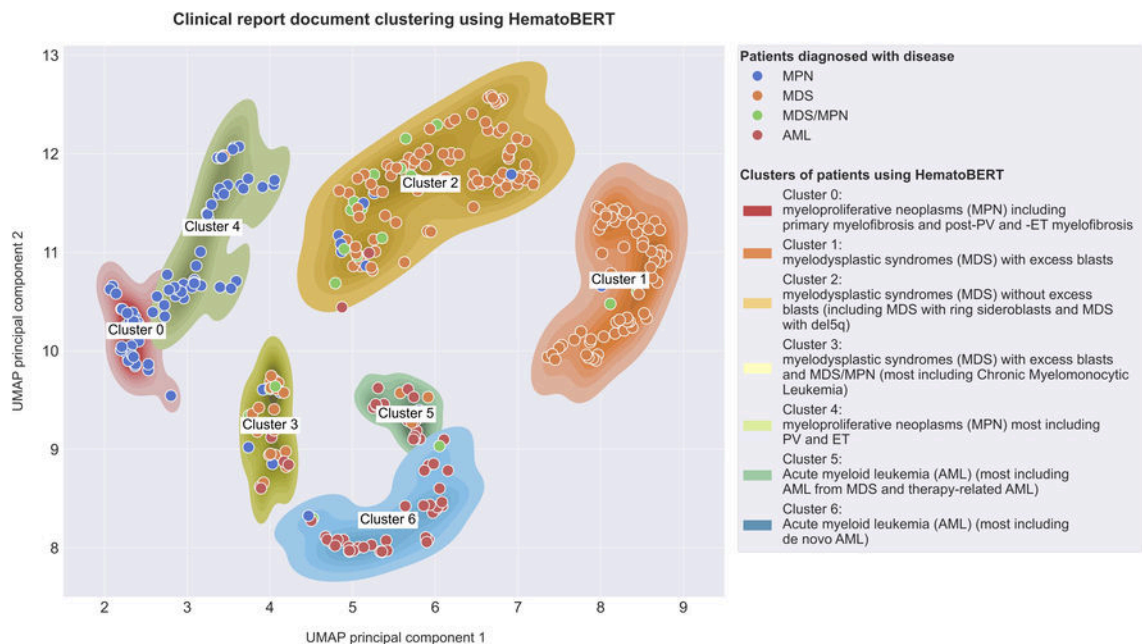
**Aims:** This project was conducted by GenoMed4all and Synthema EU consortia, with the aim to: 1) Build an AI language model specific for the hematology domain. 2) Use NLP technology to extract relevant information from clinical reports and perform unsupervised stratification of patients, in order to 3) demonstrate that the clinical report is earlier access to data relative to disease clinical phenotype and biology and provide important information for patient stratification and prediction of clinical outcomes.

**Methods:** To translate text sentences into numerical embeddings, we implemented bidirectional encoder representations from transformers (BERT) framework. To learn text representations and correlations within data, we performed domain-adaptation by fine-tuned pre-trained model on hematological clinical reports of patients with myeloproliferative neoplasms (MPN), myelodysplastic syndrome (MDS) and acute myeloid leukemia (AML). Patient stratification was performed by HDB-SCAN clustering on text embedding encoded by BERT (HematoBERT). Clusters validation was performed by assessing patients' diagnosis and survival probability. Finally, we compared domain-tuned HematoBERT vs pre-trained non-contextualized models.

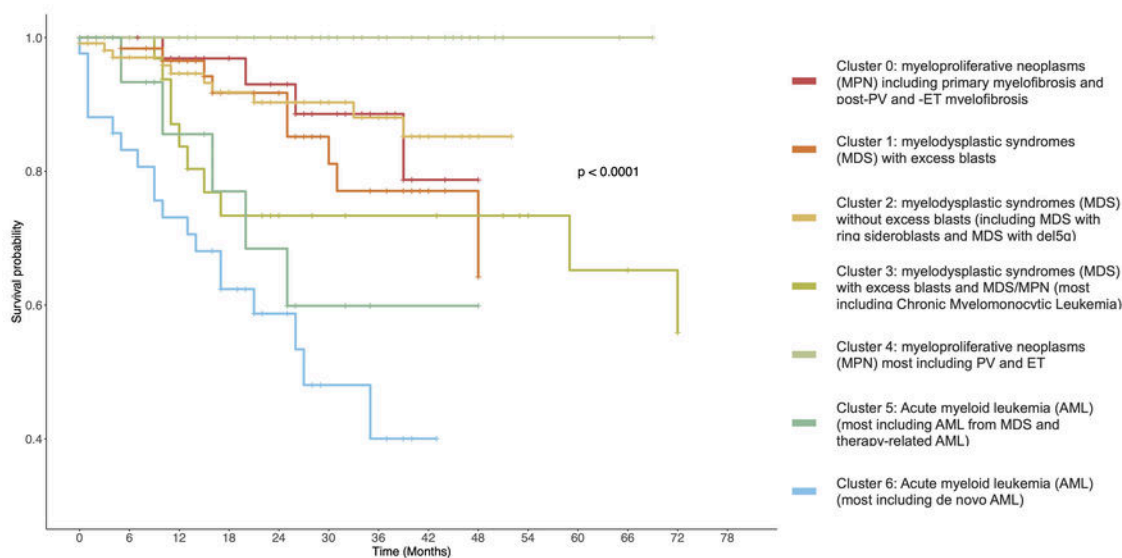
**Results:** We implemented HematoBERT based on the bert-base-multilingual-uncased version of BERT. Training data were hematological text reports of 1,328 patients. During fine-tuning, texts were tokenized, then we randomly replaced 15% of the tokens with masked tokens, training the model to predict them. We performed stratification using clinical reports from a validation cohort of 360 patients. We identified 7 clusters, defined according to similar words in meaning that were placed in a specific topic. We extracted the most important words and concepts for each cluster (topic) and we summarized them into effective descriptions for each group of patients. Two clusters included MDS patients with excess blasts, and without excess blasts with ring sideroblasts and del5q (n=69, n=115). One cluster included patients with excess blasts and MDS/MPN (n= 33). Two clusters included MPN patients with primary and secondary myelofibrosis, and MPN patients most including subjects affected with polycythemia vera and essential thrombocythemia (n=35, n=46). Two clusters included patients with AML from MDS and therapy-related AML, and patients with de novo AML (n=22, n=42). Clinical validation was performed based on the diagnosis and survival probability of patients assigned to clusters. Patients' diagnoses were compatible with the cluster assignment (Figure 1). Frequency of gene mutations (as assessed by targeted Next-Generation Sequencing) among different clusters reflected the well-known genotypic-phenotypic associations in MDS, MPN and AML. Kaplan-Maier curves indicated significative risk stratification in clusters in terms of survival probability (Figure 2), similar to stratifications performed on clinical and genomic data. Finally, we evaluate the domain adaptation by comparing the model to other pre-trained non-contextualized ones. Pseudo perplexity score (PPS), accuracy and F1 score were calculated to quantify how good the models are when they see new data, predicting the next word given the context of the sentence. HematoBERT obtained high PPS, accuracy and F1 scores, outperforming the other models also trained on generic clinical domains.

**Conclusion:** Domain-adapted language models are able to understand contexts and correlations in documents. HematoBERT can be used to extract relevant features from clinical reports. This data layer is relevant to perform disease stratification of patients based on clinical and genomic information and could be integrated into next-generation multimodal models of personalized medicine.

**Disclosures Santoro:** Amgen: Speakers Bureau; Abbvie: Speakers Bureau; Roche: Speakers Bureau; Takeda: Speakers Bureau; Merck MSD: Membership on an entity's Board of Directors or advisory committees, Speakers Bureau; Bayer: Membership on an entity's Board of Directors or advisory committees, Speakers Bureau; Eisai: Membership on an entity's Board of Directors or advisory committees, Speakers Bureau; Pfizer: Membership on an entity's Board of Directors or advisory committees, Speakers Bureau; Gilead: Membership on an entity's Board of Directors or advisory committees, Speakers Bureau; Servier: Membership on an entity's Board of Directors or advisory committees, Speakers Bureau; BMS: Membership on an entity's Board of Directors or advisory committees, Speakers Bureau; Incyte: Consultancy; Sanofi: Consultancy; Celgene (BMS): Speakers Bureau; AstraZeneca: Speakers Bureau; Eli Lilly: Speakers Bureau; Sandoz: Speakers Bureau; Novartis: Speakers Bureau; Arqule: Other. **Kordasti:** Novartis: Honoraria, Membership on an entity's Board of Directors or advisory committees; Beckman Coulter: Honoraria; MorphoSys: Research Funding. **Santini:** BMS, Abbvie, Geron, Gilead, CTI, Otsuka, servier, janssen, Syros: Membership on an entity's Board of Directors or advisory committees. **Platzbecker:** Servier: Consultancy, Honoraria, Research Funding; Geron: Consultancy, Research Funding; MDS Foundation: Membership on an entity's Board of Directors or advisory committees; Celgene: Honoraria; Merck: Research Funding; Syros: Consultancy, Honoraria, Research Funding; Fibrogen: Research Funding; Amgen: Consultancy, Research Funding; Novartis: Consultancy, Honoraria, Research Funding; Jazz: Consultancy, Honoraria, Research Funding; Takeda: Consultancy, Honoraria, Research Funding; AbbVie: Consultancy; Curis: Consultancy, Research Funding; Silence Therapeutics: Consultancy, Honoraria, Research Funding; Janssen Biotech: Consultancy, Research Funding; BMS: Research Funding; Bristol Myers Squibb: Consultancy, Honoraria, Membership on an entity's Board of Directors or advisory committees, Other: travel support; medical writing support, Research Funding; BeiGene: Research Funding; Roche: Research Funding. **Diez-Campelo:** BMS/Celgene: Consultancy, Honoraria, Membership on an entity's Board of Directors or advisory committees, Other: Advisory board fees; Novartis: Consultancy, Honoraria, Membership on an entity's Board of Directors or advisory committees; Gilead Sciences: Other: Travel expense reimbursement; GSK: Consultancy, Membership on an entity's Board of Directors or advisory committees. **Fenaux:** Janssen: Consultancy, Honoraria, Research Funding; AbbVie: Consultancy, Honoraria, Research Funding; Bristol Myers Squibb: Consultancy, Honoraria, Research Funding; Novartis: Consultancy, Honoraria, Research Funding; Jazz: Consultancy, Honoraria, Research Funding; French MDS Group: Honoraria. **Zeidan:** Shattuck Labs: Research Funding; Gilead: Consultancy, Honoraria; Celgene/BMS: Consultancy, Honoraria; AbbVie: Consultancy, Honoraria; Astex: Research Funding; Incyte: Consultancy, Honoraria; Lox Oncology: Consultancy, Honoraria; Foran: Consultancy, Research Funding; BeyondSpring: Consultancy, Honoraria; BioCryst: Consultancy, Honoraria; Notable: Consultancy, Honoraria; Kura: Consultancy, Honoraria; Tyme: Consultancy, Honoraria; Schrödinger: Consultancy, Honoraria; Zentalis: Consultancy, Honoraria; Mendus: Consultancy, Honoraria; Orum: Consultancy, Honoraria; Syndax: Consultancy, Honoraria; Epizyme: Consultancy, Honoraria; Genentech: Consultancy, Honoraria; Janssen: Consultancy, Honoraria; Amgen: Consultancy, Honoraria; Taiho: Consultancy, Honoraria; Geron: Consultancy, Honoraria; Daiichi Sankyo: Consultancy, Honoraria; Astellas: Consultancy, Honoraria; Novartis: Consultancy, Honoraria; Boehringer-Ingelheim: Consultancy, Honoraria; Servier: Consultancy, Honoraria; Agios: Consultancy, Honoraria; Pfizer: Consultancy, Honoraria; Seattle Genetics: Consultancy, Honoraria; Ionis: Consultancy, Honoraria; Takeda: Consultancy, Honoraria; Otsuka: Consultancy, Honoraria; Chiesi: Consultancy, Honoraria; ALX Oncology: Consultancy, Honoraria; Regeneron: Consultancy, Honoraria; Jazz: Consultancy, Honoraria; Syros: Consultancy, Honoraria. **Haferlach:** MLL Munich Leukemia Laboratory: Current Employment, Other: Equity Ownership. **Della Porta:** Bristol Myers Squibb: Honoraria, Membership on an entity's Board of Directors or advisory committees.



**Figure 1. Clinical reports-based clustering of patients with MDS, MPN and AML using HematoBERT.** UMAP 2-dimensional text embedding encoded by fine-tuned HematoBERT (each dot represents a patient) whose location is defined on the basis of the textual representation and correlation. The figure shows the clustering of patients' (noise cluster was excluded) clinical records according to similar words and concepts. Stratification of patients by diagnosis is consistent with cluster groupings.



**Figure 2. Kaplan-Meier probability estimates of Overall Survival (OS) stratified by clusters identified using HematoBERT.**

**Figure 1**

<https://doi.org/10.1182/blood-2023-188292>