

# Diversity Minute

But first, an announcement...

**March 4, 2022**



Dear Haihao,

Congratulations! I am pleased to inform you that you have been admitted to the Master of Education program in Learning Design, Innovation, and Technology at the Harvard Graduate School of Education to pursue full-time study for the 2022-2023 academic year.

You will be a part of a cohort of students who bring with them not only impressive professional experience and excellent academic training, but also dedication to the profession of education. You are to be commended on the fine record of achievement by which you earned your place in the class.

# About me



- From Shenzhen, China
  - and Stamford, CT, and Hong Kong
- Rice University
  - BS in Materials Science and NanoEngineering
  - BA in Mathematics
- Diverse research experiences abroad in undergrad



Climbing Mount Fuji



Skiing in the French Alps



Tomb of Confucius

# Information Extraction from Materials Science Literature: Machine Learning Tasks and Methods

Haihao Liu

March 7, 2022



# Materials Informatics

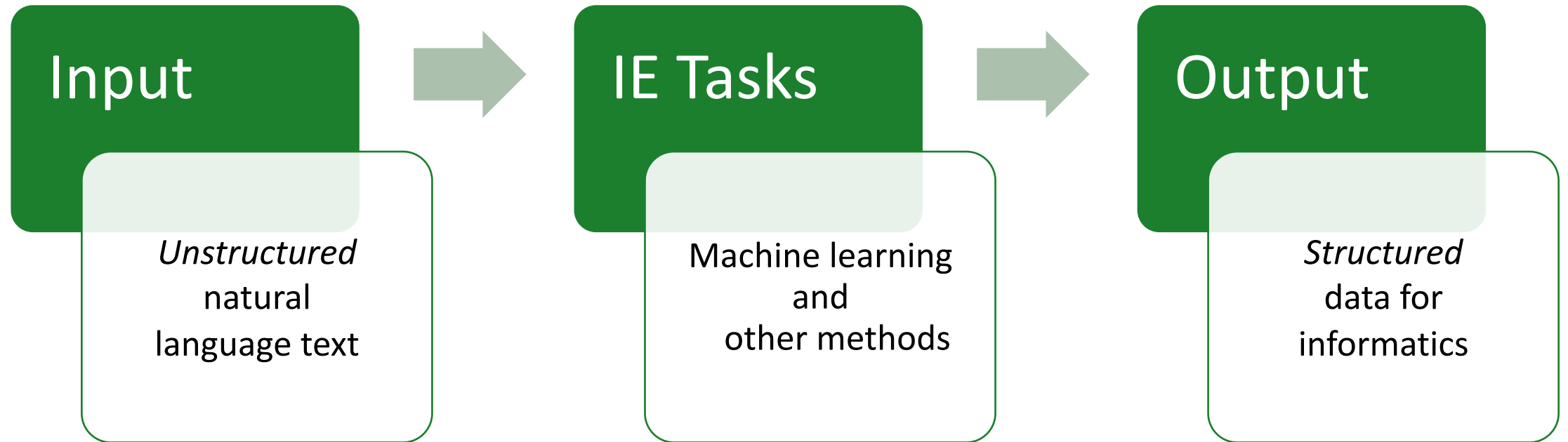
## Motivation

- “Fourth paradigm” after experiments, theory, simulation
- Need high-quality structured data for materials informatics
- Sparsity and scarcity key issues



# Information Extraction

Background



# Natural Language Processing

- Empirical linguistics since 90s
- Applied fields – chemistry, polymers, biology/medical
- Contextualized word embeddings (ELMo, BERT)



## Information Extraction from Materials Science Literature

Task 1:

Sample name  
recognition

Task 2:

Aluminum alloy  
data extraction

Task 3:

Key phrase ranking  
and analysis

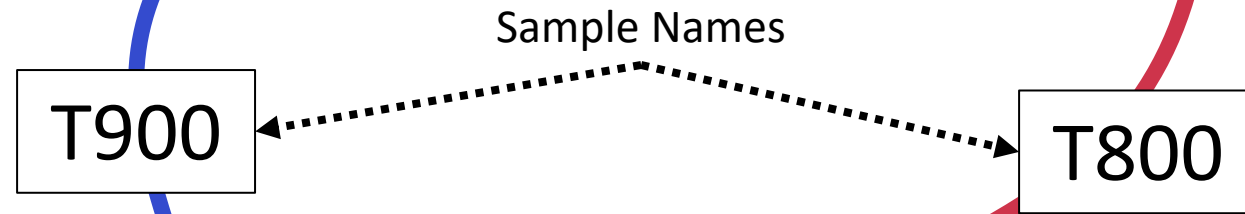
# Task 1:

## Sample name recognition

# Task Definition

## Task 1

Two low-density steel specimens were prepared for experiments by cold-rolling and annealing at 800 °C (T800) and 900 °C (T900) for 2 min in an infrared heating furnace. The macroscopic tensile behaviors of both



scale tensile behavior was observed in both steels. The yield strength of T800 is 718 MPa, which is much higher than the 561 MPa obtained for T900. While T800 shows

# Hyperparameter Tuning

## Task 1

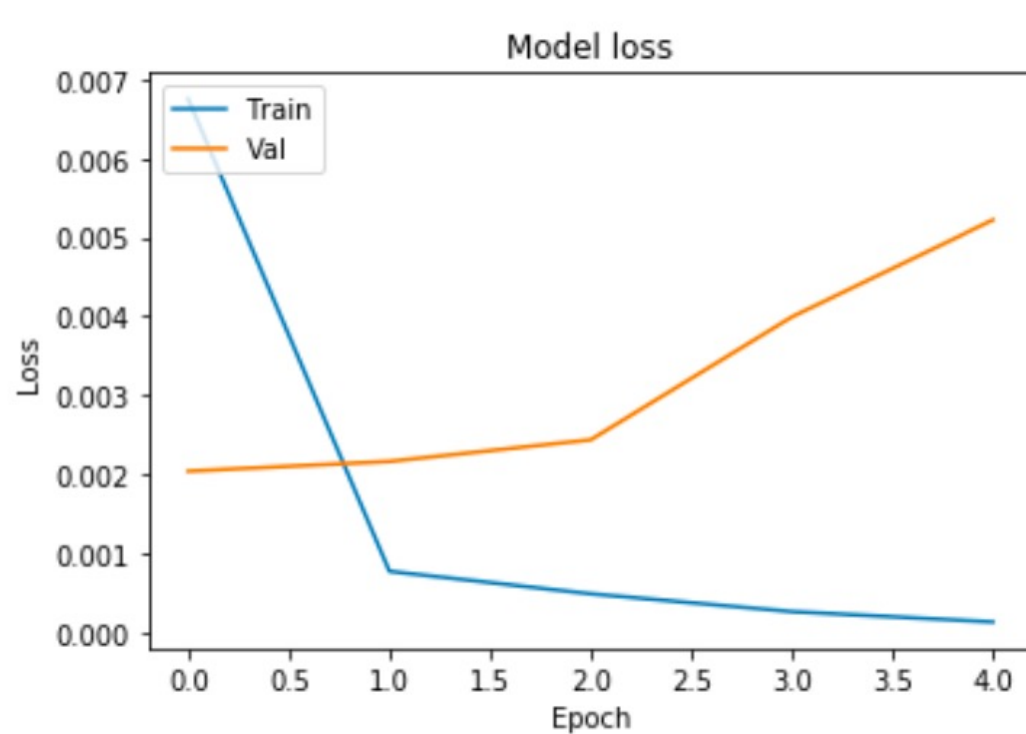
$$\text{FL}(y, \hat{y}) = -\alpha \sum_{c=1}^m (1 - \tilde{y}_c)^\gamma y_c \log(\tilde{y}_c) \in \mathbb{R}^n$$

$\gamma \setminus lr$	$10^{-7}$	$10^{-6}$	$10^{-5}$	$10^{-4}$	$10^{-3}$
0.0	0.0536	<b>0.4857</b>	0.5035	0.4463	0.0
2.0	0.0929	0.4739	<b>0.5048</b>	<b>0.4770</b>	0.0
4.0	<b>0.2045</b>	0.4795	0.4926	0.4388	0.0

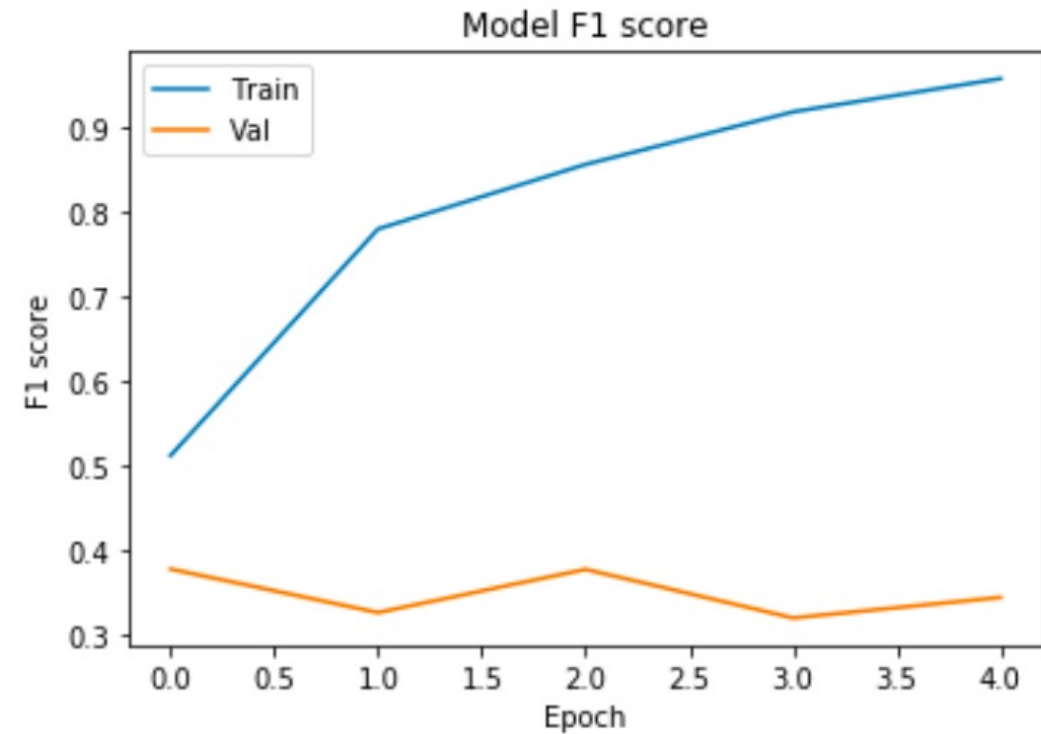
Table 3: F1 results of hyperparameter grid search

# Hyperparameter Tuning

## Task 1



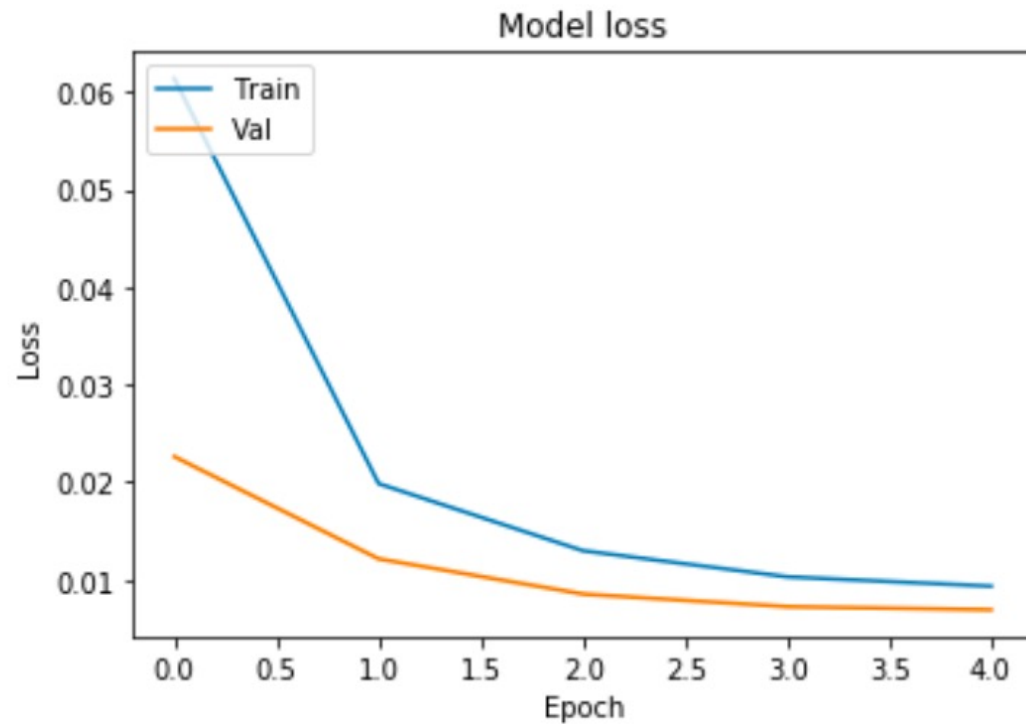
(a) Focal loss over time ( $lr = 10^{-4}, \gamma = 2.0$ )



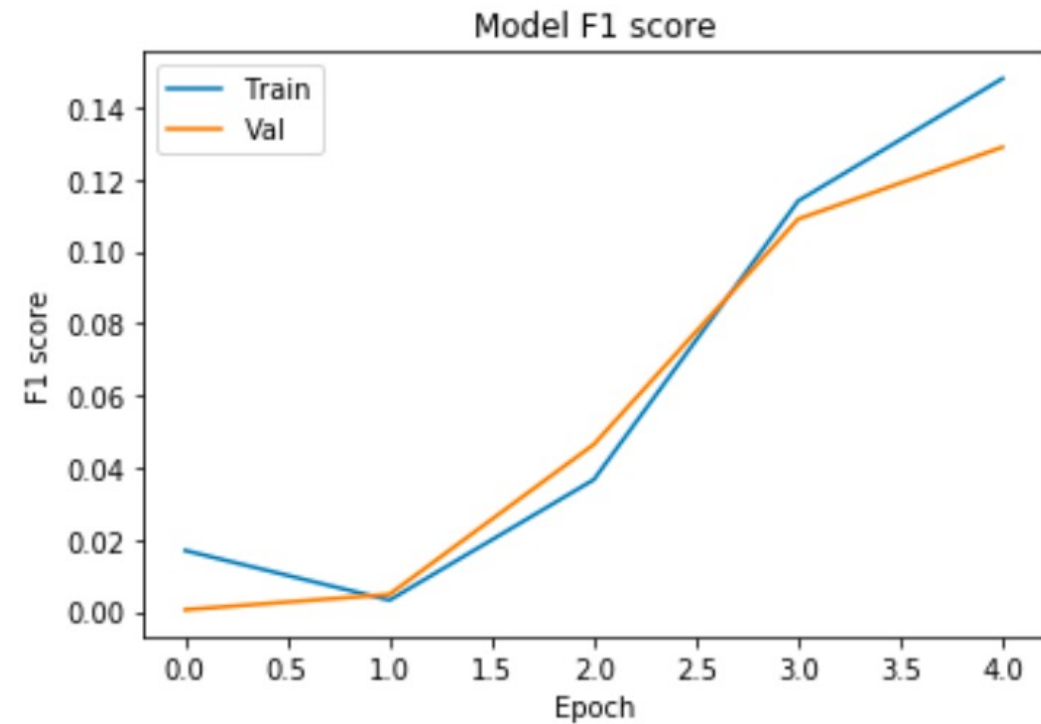
(b) F1 over time ( $lr = 10^{-4}, \gamma = 2.0$ )

# Hyperparameter Tuning

## Task 1



(c) Focal loss over time ( $lr = 10^{-7}, \gamma = 4.0$ )



(d) F1 over time ( $lr = 10^{-7}, \gamma = 4.0$ )

# Model Interpretation

**y=B** (probability **0.895**, score **3.429**) top features

Contribution?	Feature
+13.798	Highlighted in text (sum)
-10.369	<BIAS>

[CLS] The affinity follows the sequence **H ##60 ##0 - 0 . 5 N >** H ##60 ##0 - 3 N > H ##60 ##0 - 5 N > H ##60 ##0 . [SEP]

**y=B** (probability **0.959**, score **4.469**) top features

Contribution?	Feature
+14.062	Highlighted in text (sum)
-9.594	<BIAS>

[CLS] All of the data were collected by adding 10 mg solid material into 10 m ##L H ##A solution , which was then kept at room temperature for 6 h . A few points worth noting are : ( 1 ) comparing Re value of **PA ##C - C** and PA ##C - P , it is shown that PA ##C - C exhibits a slight better performance than PA ##C - P . [SEP]



# Model Interpretation

## Task 1

**y=O** (probability **0.555**, score **0.599**) top features

Contribution?	Feature
+2.194	<BIAS>
-1.595	Highlighted in text (sum)

[CLS] After the four consecutive cat ##alytic runs , we observed a higher decrease of conversion for Amber ##ly ##st - 15 than for S ##BA - 15 - M ##w ##S and SM ##w - AG catalyst ##s , which showed just a slight conversion decrease ( Fi ##g . [SEP]

**y=B** (probability **0.441**, score **0.052**) top features

Contribution?	Feature
+4.392	Highlighted in text (sum)
-4.340	<BIAS>

[CLS] After the four consecutive cat ##alytic runs , we observed a higher decrease of conversion for Amber ##ly ##st - 15 than for S ##BA - 15 - M ##w ##S and SM ##w - AG catalyst ##s , which showed just a slight conversion decrease ( Fi ##g . [SEP]

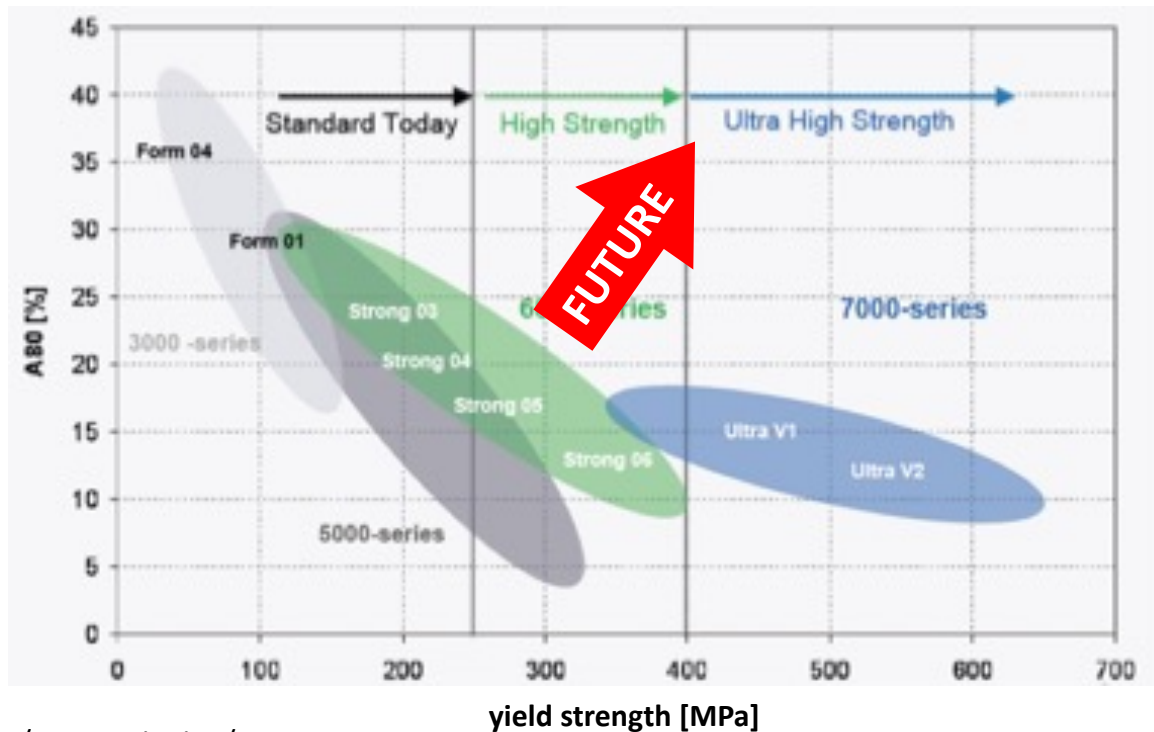
# Task 2:

## Aluminum alloy data extraction

# Aluminum alloy design

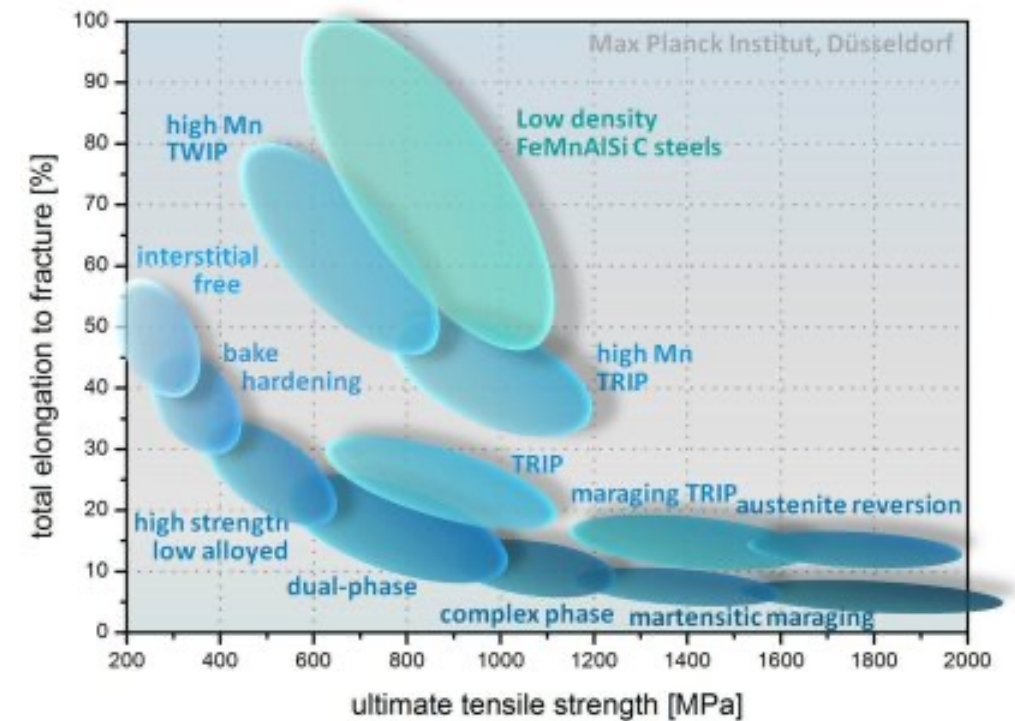
- Second-most produced metal after steel, more room to innovate

## Aluminum alloys



TWIP/TRIP: Twinning/  
transformation-induced plasticity

## Steel alloys



# Task Definition

## Task 2

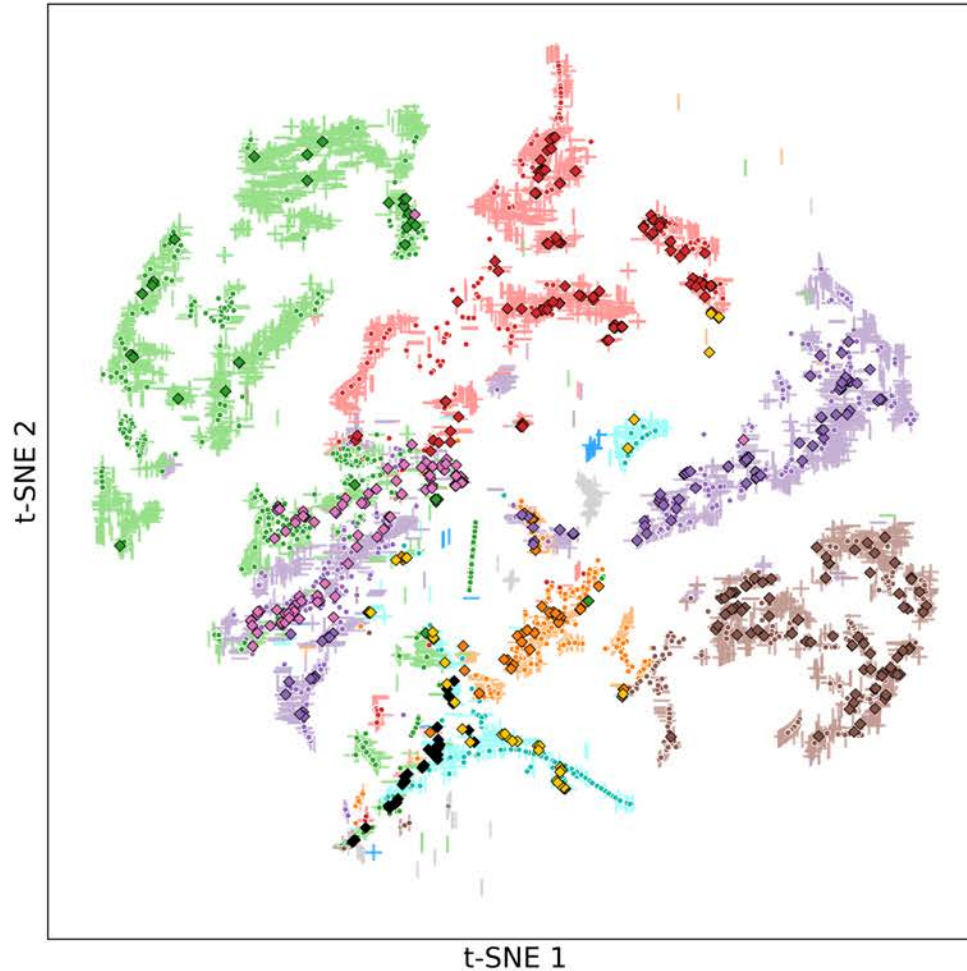
**Table 1 Alloying elements in wrought aluminum and aluminum alloys**

Representative list of common wrought aluminum alloys

Aluminum Association designation	Unified Numbering System (UNS) designation	Composition, maximum unless a range or minimum is specified(a), wt%							Al, min or bal
		Si	Fe	Cu	Mn	Mg	Zn	Other specified alloying elements	
1050	A91050	0.25	0.40	0.05	0.05	0.05	0.05	...	99.50
1060	A91060	0.25	0.35	0.05	0.03	0.03	0.05	...	99.60
1145	A91145	0.55 Si + Fe		0.05	0.05	0.05	0.05	...	99.45
1175	A91175	0.15 Si + Fe		0.10	0.02	0.02	0.04	...	99.75
1200	A91200	1.00 (Si + Fe)		0.05	0.05	...	0.10	...	99.0
1230	A91230	0.70 Si + Fe		0.10	0.05	0.05	0.10	...	99.30
1235	A91235	0.65 Si + Fe		0.05	0.05	0.05	0.10	...	99.35
1345	A91345	0.30	0.40	0.10	0.05	0.05	0.05	...	99.45
1350	A91350	0.10	0.40	0.05	0.01	...	...	...	99.50
2011	A92011	0.40	0.7	5.0–6.0	...	...	0.30	0.20–0.6% Bi; 0.20–0.6% Pb	bal
2014	A92014	0.50–1.2	0.7	3.9–5.0	0.40–1.2	0.20–0.8	0.25	...	bal
2017	A92017	0.20–0.8	0.7	3.5–4.5	0.40–1.0	0.40–0.8	0.25	...	bal
2018	A92018	0.9	1.0	3.5–4.5	0.20	0.45–0.9	0.25	1.7–2.3Ni	bal
2024	A92024	0.50	0.50	3.8–4.9	0.30–0.9	1.2–1.8	0.25	...	bal
2025	A92025	0.50–1.2	1.0	3.9–5.0	0.40–1.2	0.05	0.25	...	bal
2036	A92036	0.50	0.50	2.2–3.0	0.10–0.40	0.30–0.6	0.25	...	bal
2117	A92117	0.8	0.7	2.2–3.0	0.20	0.20–0.50	0.25	...	bal
2124	A92124	0.20	0.30	3.8–4.9	0.30–0.9	1.2–1.8	0.25	...	bal
2218	A92218	0.9	1.0	3.5–4.5	0.20	1.2–1.8	0.25	1.7–2.3Ni	bal

# Alloy Compositions

## Task 2

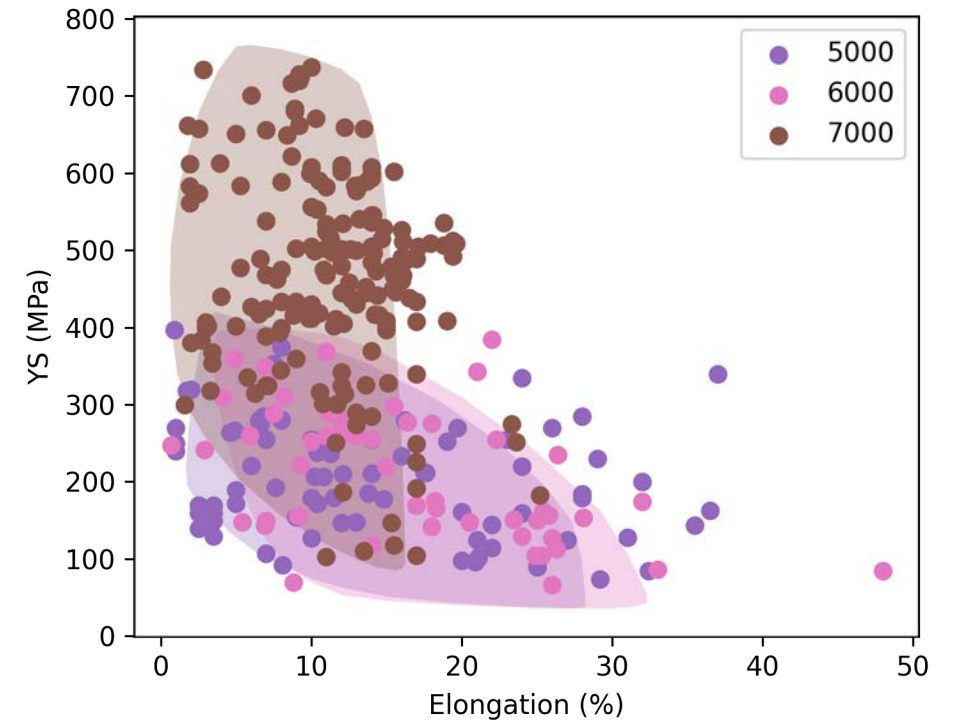
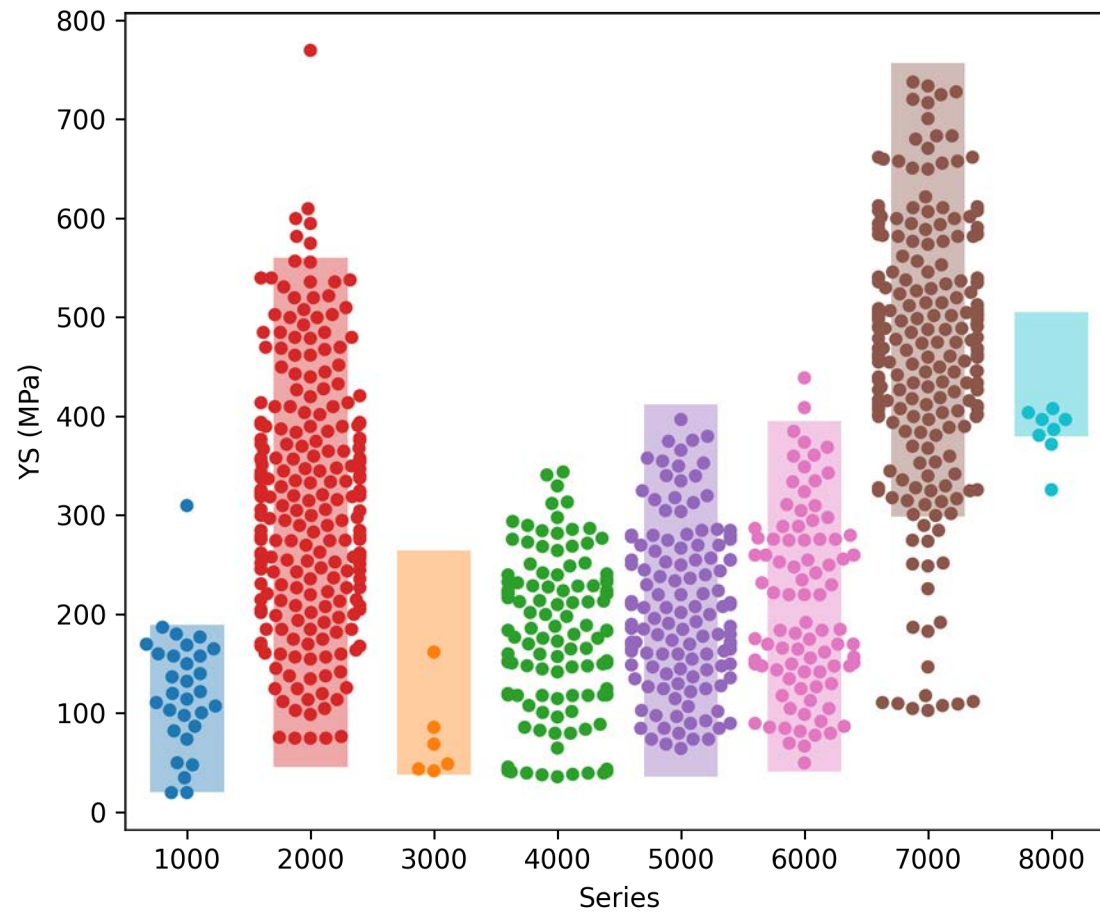


Registered Alloys	Journal Texts	Journal Tables	Patents
◆ 1000	—	—	—
◆ 2000	Cu	— Cu	• Cu
◆ 3000	Mn	— Mn	• Mn
◆ 4000	Si	— Si	• Si
◆ 5000	Mg	— Mg	• Mg
◆ 6000	—	—	—
◆ 7000	Zn	— Zn	• Zn
◆ 8000	—	—	—
	Cr	— Cr	—
	Fe	— Fe	• Fe
	Ti	— Ti	• Ti

n = 14,884

# Alloy Properties

## Task 2



n = 1,278



Task 3:

Key phrase ranking and analysis



# Top key phrases per series

5000	6000	7000
friction stir welding	t6 aluminium alloy	aluminum alloy 7075
5083 aluminum alloy	severe plastic deformation	7050 aluminum alloy
5052 aluminum alloy	finite element method	stress corrosion cracking
ultimate tensile strength	friction stir processing	fatigue crack growth
aluminum alloy 5083	tool rotational speed	ultimate tensile strength
friction stir processing	finite element analysis	aluminium alloy 7075
aluminum alloy 5052	energy dispersive x	aluminum alloy 7050
heat affected zone	6082 aluminum alloy	average grain size
resistance spot welding	experimental results show	heat affected zone
finite element analysis	differential scanning calorimetry	t7451 aluminum alloy
aluminum alloy sheet	response surface methodology	t6 aluminium alloy
aluminium alloy 5083	6061 aluminum alloys	finite element method
aluminum alloy sheets	fatigue crack growth	7075 aluminum alloys
5754 aluminum alloy	metal matrix composites	experimental results show
average grain size	aa6061 aluminum alloy	solution heat treatment
tool rotational speed	average grain size	fatigue crack initiation
stress corrosion cracking	friction stir welded	7085 aluminum alloy
finite element method	6061 al alloy	crack growth rate
strain rate sensitivity	finite element model	high strength aluminum alloys
5182 aluminum alloy	solution heat treatment	strain rate range
energy dispersive x	aluminum alloy 6063	energy dispersive x
digital image correlation	friction stir spot welding	7075 al alloy

# Top studied alloys by series

## Task 3

2024	977.655614	3003	100.689273				
2219	201.870831	3004	29.207618	6061	1262.18073	7075	1040.16121
2014	87.982418	3104	24.944444	6063	334.403977	7050	289.886416
2124	41.268038	3105	7.75	6082	227.730373	7055	69.186334
2017	37.173993			6016	81.378912	7475	64.352704
2618	32.9	5083	427.388609	6060	58.515152	7150	56.630952
2524	27.848851	5052	316.932317	6111	55.923308	7085	49.278571
2050	11	5754	149.311418	6013	35.199856	7020	41.374747
2195	9.766667	5182	101.161663	6022	23.55	7010	40.274359
2519	9.474747	5056	20.816667	6351	16.725	7005	30.106061
2624	8	5086	16.866667	6005	15.492308	7136	16
2011	7.333333	5005	15.287879	6056	14.583333	7175	12.271429
2099	5.8	5356	11	6101	11.119048	7449	9.733333
2139	5	5456	7.166667			7003	9.383333

Scores are sums of normalized occurrences

# Top concepts specific to 7000 series

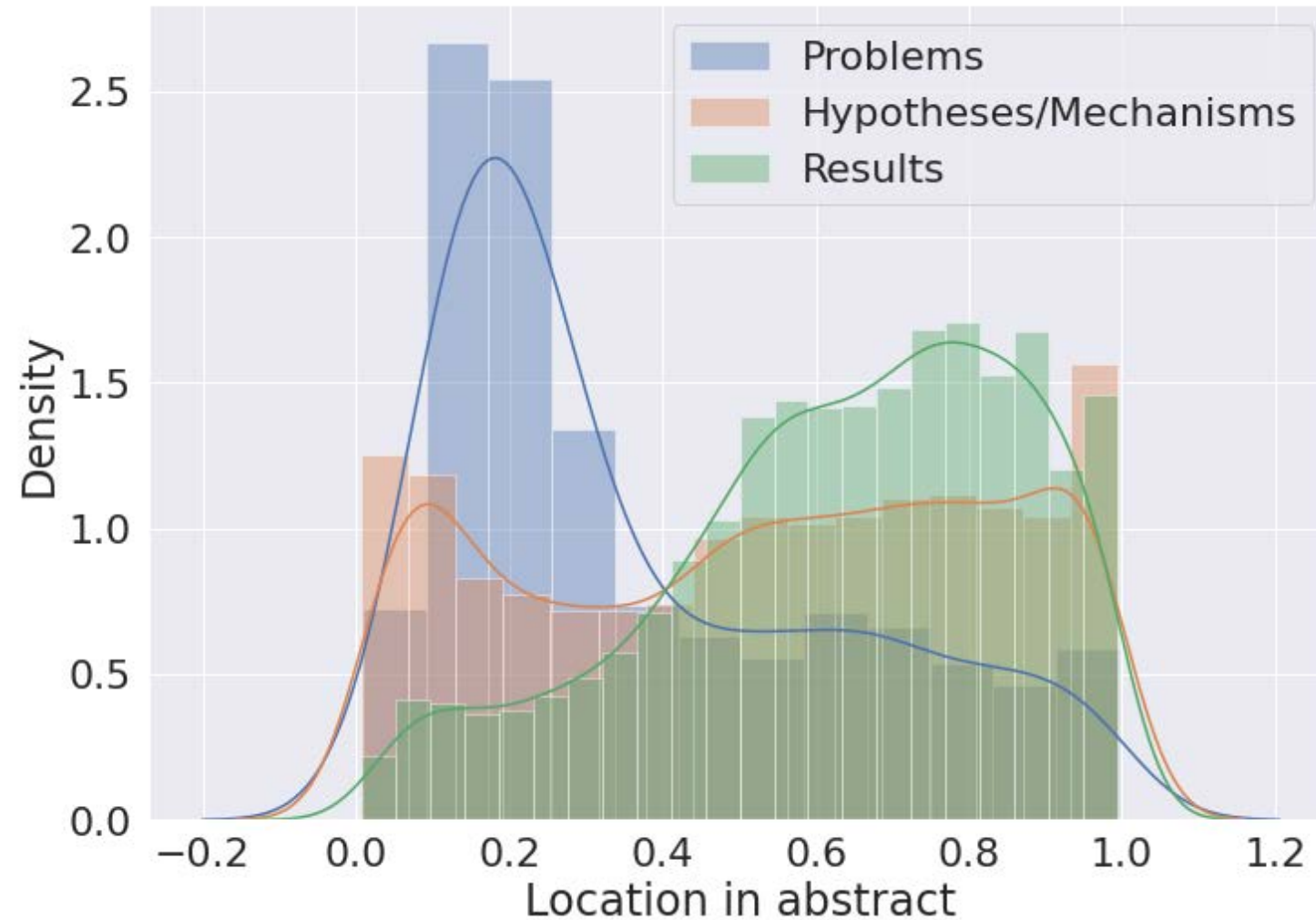
- Tempers
  - T6, T651, T7351, T7451, T76
  - cf. non-heat-treatable series e.g. 5000: H111, H116, H34, H321
- Temperatures
  - 120 °C, 470 °C, 480 °C
- Topics of study
  - stress corrosion cracking
  - fatigue crack growth
  - $\eta'$  phase
  - slow strain rate test
  - residual stress
  - creep age forming
  - high strength aluminum alloys
  - surface roughness
  - grain boundary precipitates
  - residual stress distribution

# Links between keywords

channel angular pressing	average grain size	62
laser shock peening	compressive residual stress	45
channel angular pressing	fine grain size	41
initial strain rate	grain boundary sliding	40
equal channel angular pressing	average grain size	32
channel angular pressing	initial strain rate	26
resistance spot welding	tensile shear strength	26
grain boundary sliding	high strain rate superplasticity	26
thermal expansion coefficient	metal matrix composite	26
cast aluminum alloys	dendrite arm spacing	26
initial strain rate	average grain size	23
high specific strength	metal matrix composite	23
channel angular pressing	grain boundary sliding	21
maximum tensile strength	tool rotational speed	21
response surface methodology	metal matrix composite	21
	tool rotational speed	21
channel angular pressing	high strain rate superplasticity	21
compressive residual stress	fatigue crack growth rate	20
continuous dynamic recrystallization	average grain size	19
secondary dendrite arm spacing	cast aluminum alloys	19
channel angular pressing	high angle grain boundaries	18
	continuous dynamic recrystallization	18
fatigue crack growth rate	high strength al	18

# Inferring relationships between links

## Task 3



# Conclusions and Outlook

- Three information extraction tasks
  - Sample name recognition, Al alloy data extraction, key phrase analysis
- NLP still young in its application to mat sci, many challenges remain
- To make most use of data, need to incorporate domain knowledge



# Acknowledgments





Questions?