

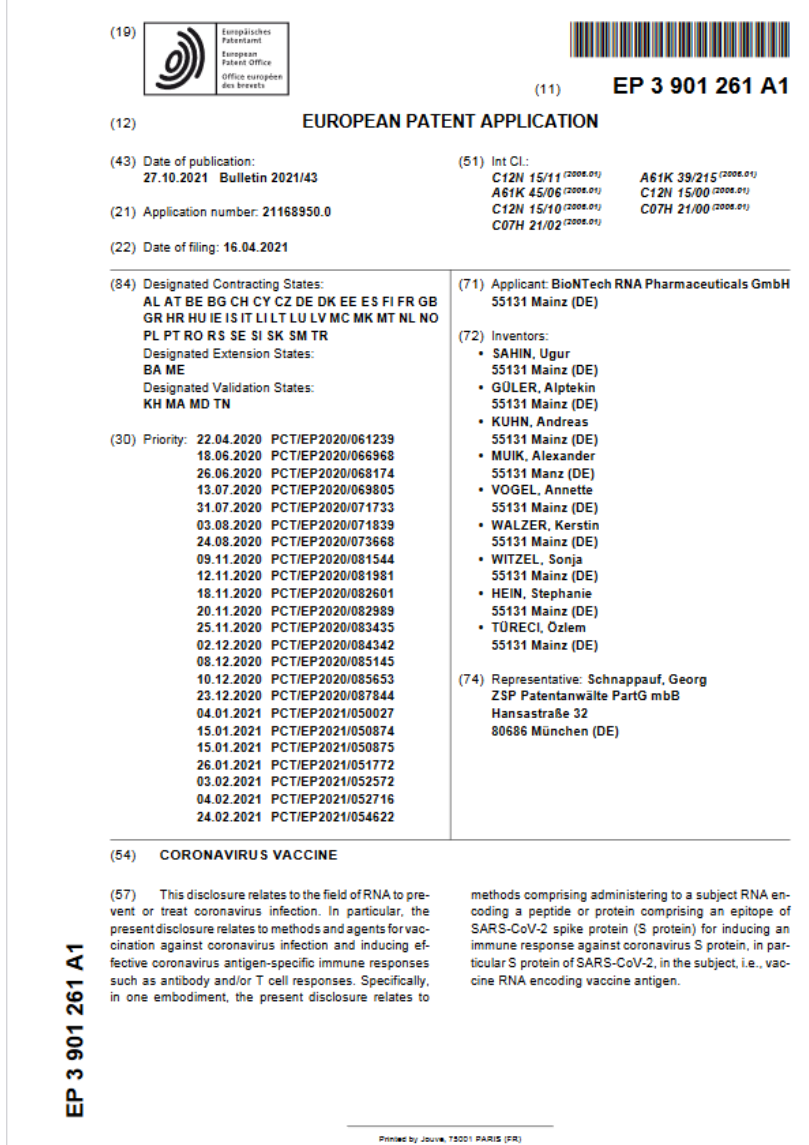
Moving On – Investigating Inventors' Ethnic Origins Using Supervised Learning

The Problem

- Patent data provides rich information about technical inventions and is frequently used in social sciences research.
- Unfortunately, it does **not disclose inventors' ethnic origins**.
- Current approaches use costly commercial tools to map inventors to ethnicities.
- Particular challenge for research at the intersection of immigration and innovation topics

Contributions

- Construct a new dataset of names that is representative for inventor origin analysis
- Train and validate an algorithm that predicts ethnic origins based on names
- Provide estimates for inventor ethnic origin compositions over time, across countries and technological areas.



Dataset: Olympic Athletes' Names

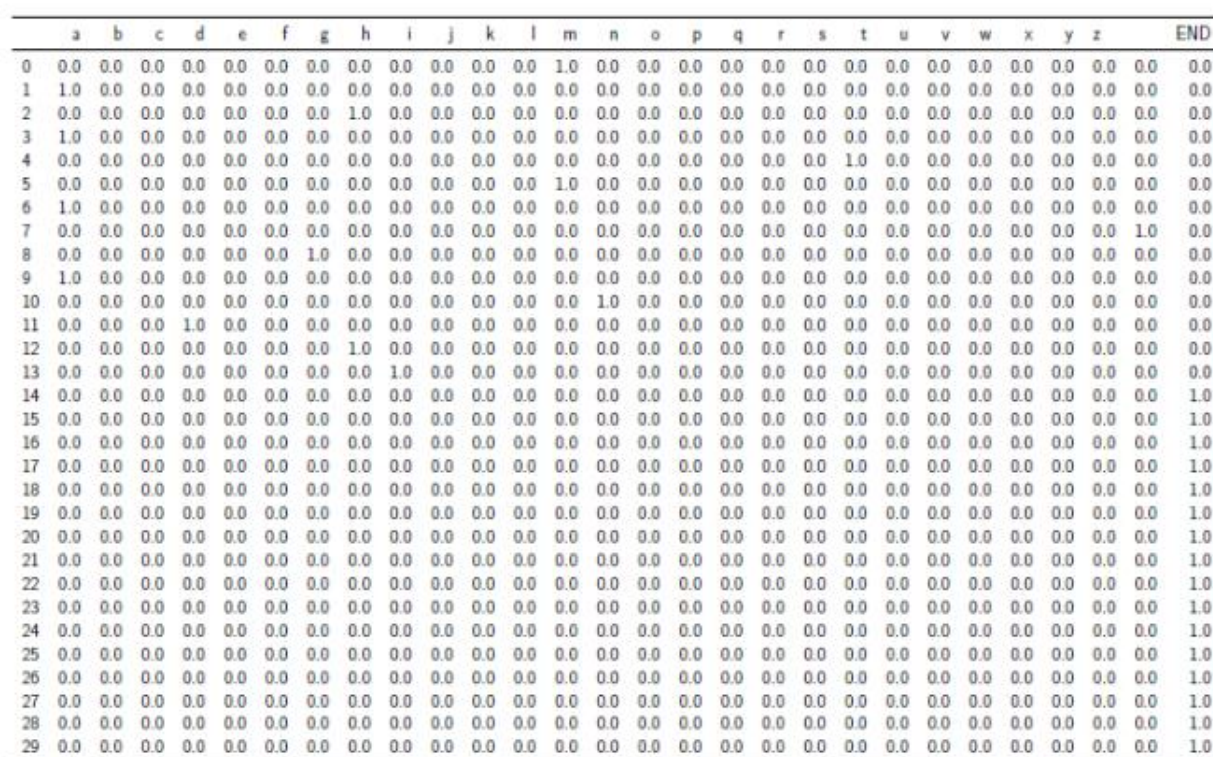
- 135'571 full names of athletes who have been starters for 191 different national teams in Olympic games from Athens 1896 to Rio 2016.
 - Label ethnic origins of these names by mapping national teams to ethnic origins.
 - Clean and enrich with inventor names to a final dataset of 95'202 labelled names.
- Use the letters a name consists of to predict the name's ethnic origin with an LSTM
- 91% F1-score across 17 ethnic origins

Table 1: Taxonomy of Ethnic Origins and National Teams

Ethnic Origin	National Teams / Countries
Anglo-Saxon	Great Britain, Ireland
Chinese	China
French	France
German	Germany
Hispanic-Iberian	Spain, Portugal, Mexico
India	India
Italian	Italy
Japanese	Japan
Korean	Korea
Arabic	Egypt, Syria, Saudi Arabia, Jordan, UAE, Tunisia, Algeria, Morocco
Persian	Iran
Slavic-Russian	Russia, Ukraine, Belarus
East-Europe	Poland, Czechoslovakia, Hungary
Balkans	Serbia, Croatia, Yugoslavia
Scandinavian	Sweden, Norway, Finland, Denmark, Iceland
South-East Asia	Vietnam, Thailand, Malaysia, Indonesia, Laos, Cambodia
Turkish	Turkey

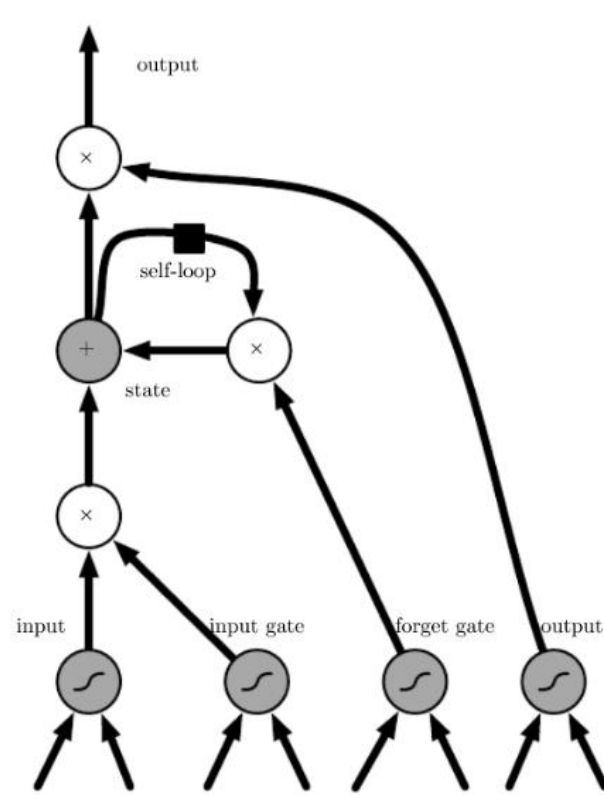
Name Transformation

Figure D.1: Example of the Encoded Name Mahatma Gandhi



LSTM Network

Figure D.2: Visualization of an LSTM cell (taken from Goodfellow et al. [2016])



Performance

Table 2: Performance of the LSTM Classification Model

Ethnic Origin	Precision	Recall	F1 Score
Overall (weighted)	0.910	0.910	0.910
Anglo-Saxon	0.859	0.880	0.873
Arabic	0.911	0.927	0.919
Balkans	0.816	0.753	0.783
Chinese	0.932	0.938	0.938
East-Europe	0.908	0.913	0.910
French	0.911	0.893	0.902
German	0.805	0.860	0.832
Hispanic-Iberian	0.908	0.925	0.916
India	0.910	0.852	0.880
Italian	0.955	0.901	0.927
Japanese	0.972	0.993	0.982
Korean	0.925	0.955	0.940
Persian	0.897	0.902	0.900
Scandinavian	0.919	0.895	0.907
Slavic-Russian	0.961	0.956	0.958
South-East Asia	0.871	0.758	0.810
Turkish	0.923	0.966	0.944

Inventors' Ethnic Origins

- Predict the ethnic origins of 2.68 million patent inventors from 1980 to 2015
- Global ethnic origin distribution has become more diverse. Mostly due to more Asian origin inventors (Chinese, Indian)
- Prevalence of foreign-origin inventors is especially high in the USA but has also increased in other high-income economies.
- Non-Western origin inventors are to be found mostly in the USA but not to European countries.

The USA has an extraordinary ability to attract global talents

vs.

Continental European countries are lagging behind

Figure 1: The Global Ethnic Origin Composition of Inventors (1980-2015)

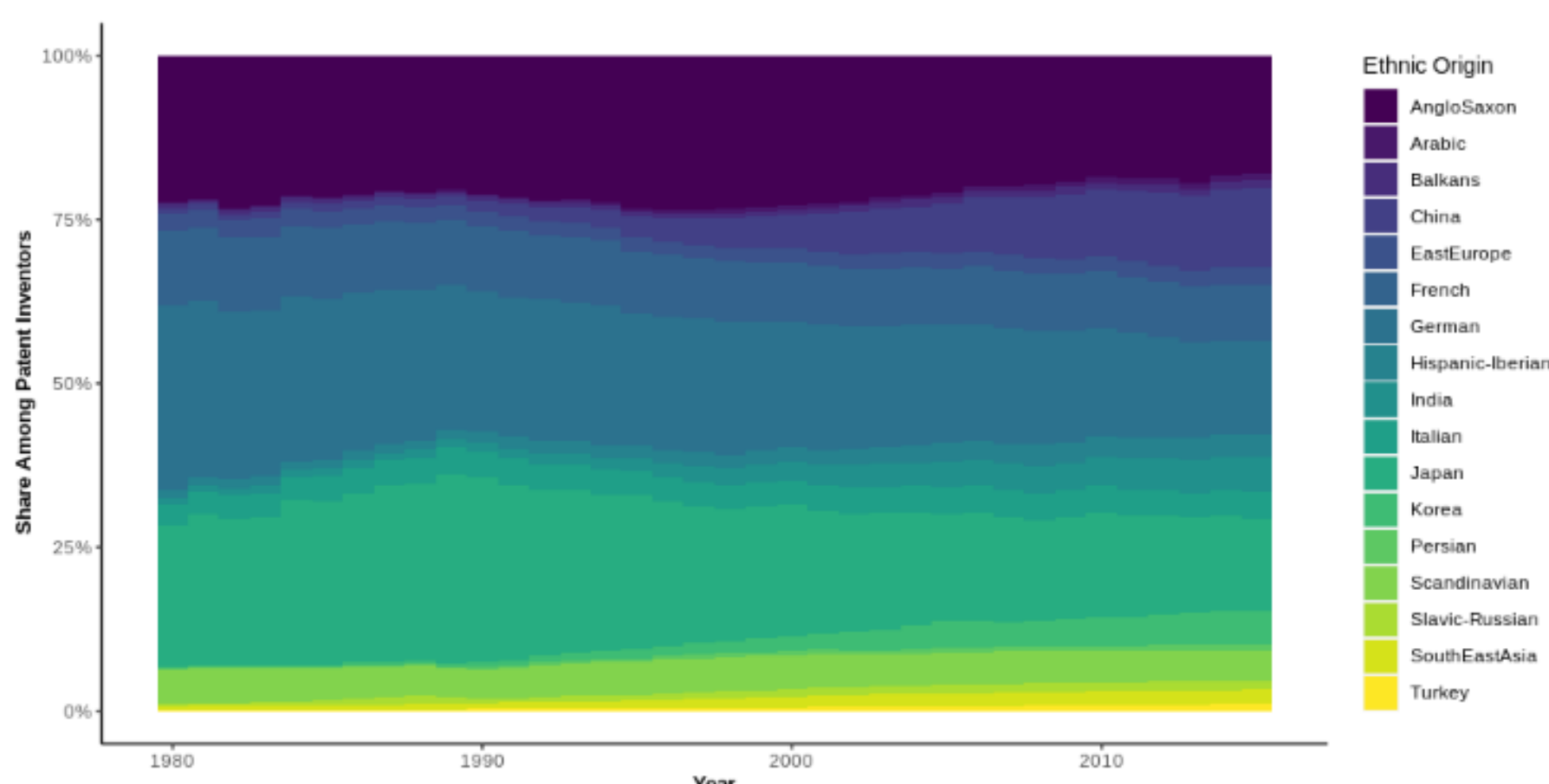
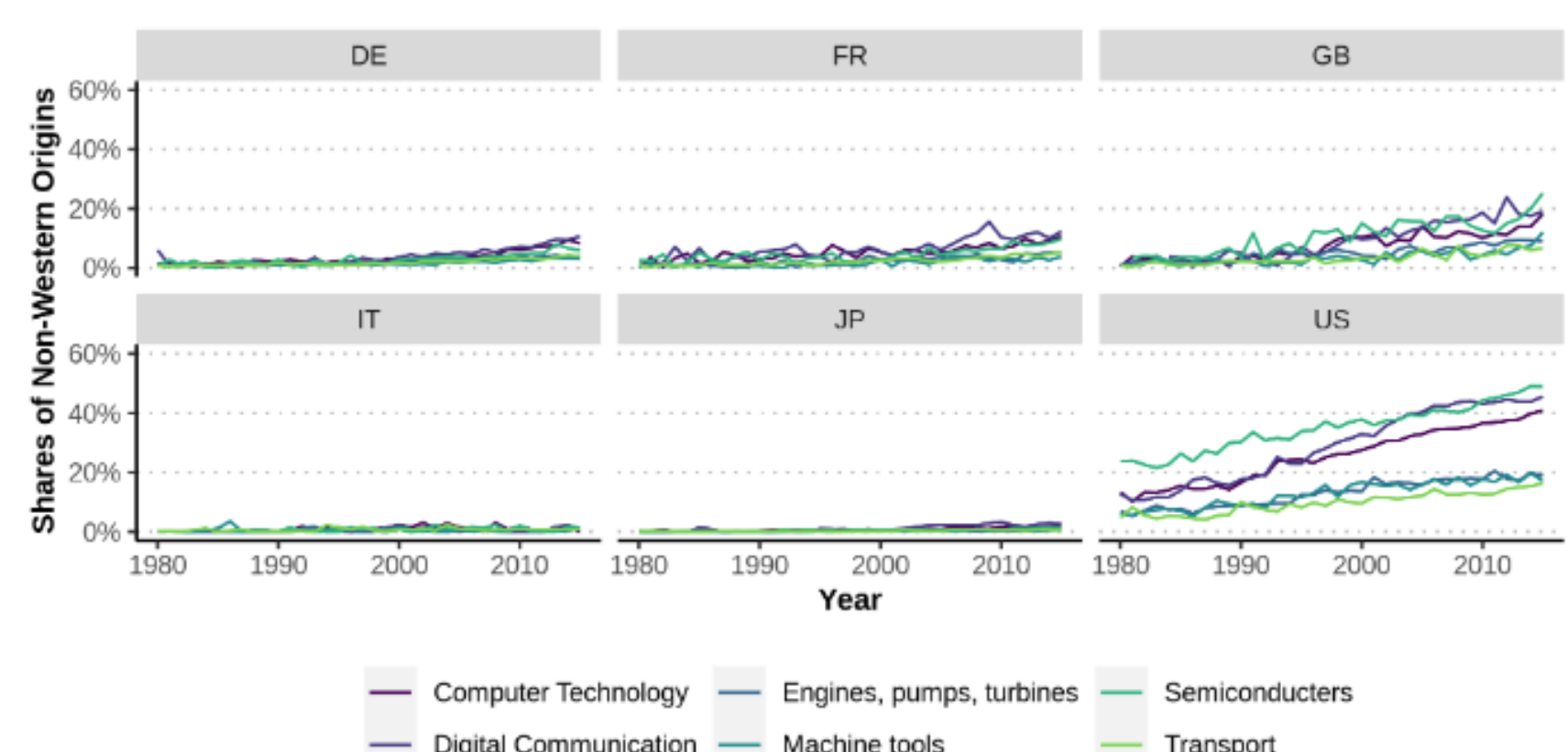


Figure 3: Prevalence of Non-Western Ethnic Origin Across Technologies (1980-2015)



Notes: The graph shows the evolution of the prevalence of 17 ethnic origins among the annual stock of patent inventors between 1980 and 2015. Formally, it plots $\pi_{j,t}^k$. The data for this plot is from the OECD and the USPTO.

Notes: The graph shows the evolution of the aggregate prevalence of non-western ethnic origins among the annual stock of patent inventors in six technology fields for six high-income economies. Formally, it plots $\sum_{k=1}^K \pi_{j,t}^k$, for each technology field j , with k corresponding to the following non-western ethnic origins: Arabic, Chinese, Indian, Persian, Slavic-Russian, Turkish and South-East Asian. The data for this plot is from the OECD and the USPTO.