

PREDICTING INNOVATIVE FIRMS USING WEB MINING AND DEEP LEARNING

Jan Kinne & David Lenz

May 2019

IGL 2019



MOTIVATION: SHORTCOMINGS OF TRADITIONAL INNOVATION INDICATORS

- ❖ Timeliness
- ❖ Coverage
- ❖ Data collection costs
- ❖ Granularity

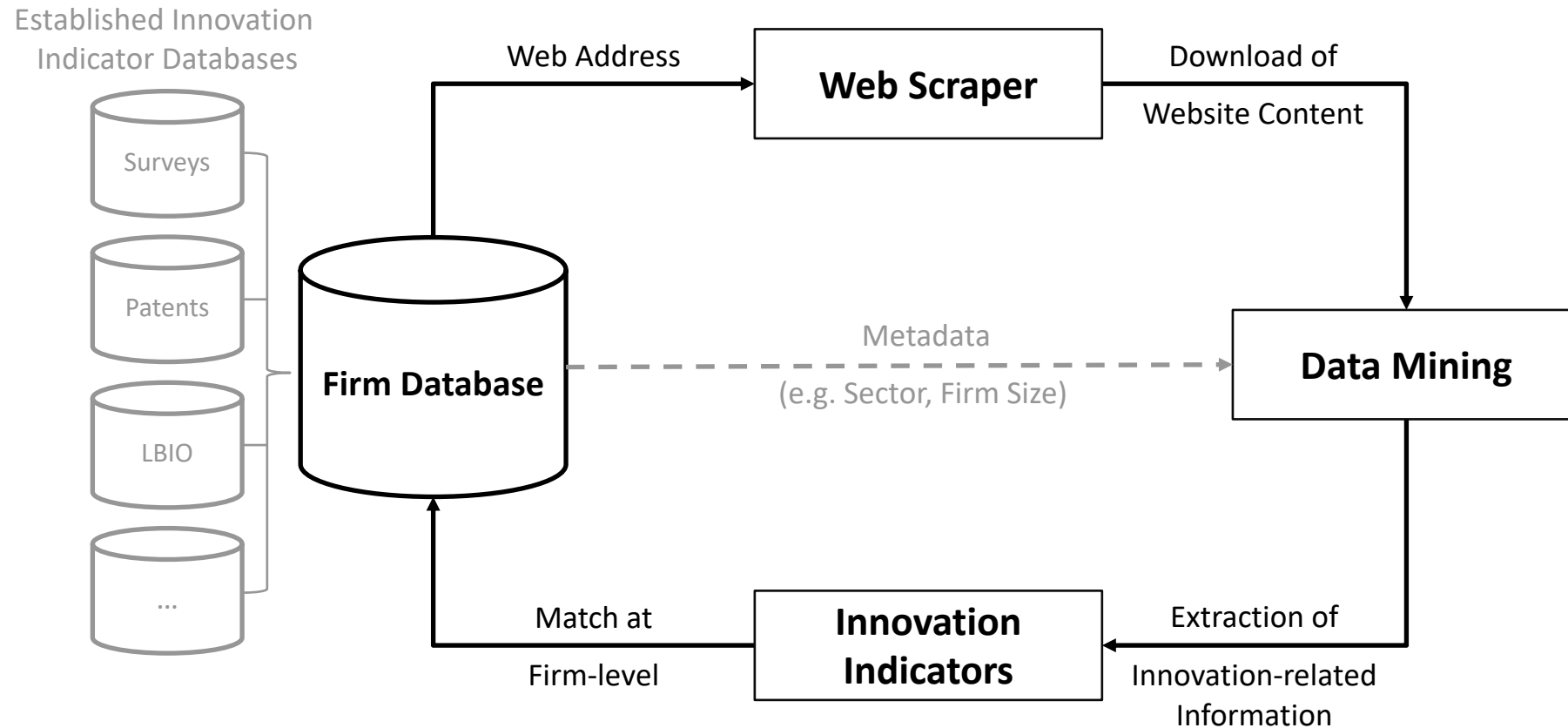
STI policy making requires an accurate and timely picture of the current state of the STI system in order to plan and evaluate policy measures in an evidence-based manner.

MOTIVATION: WEB MINING OF FIRM WEBSITES

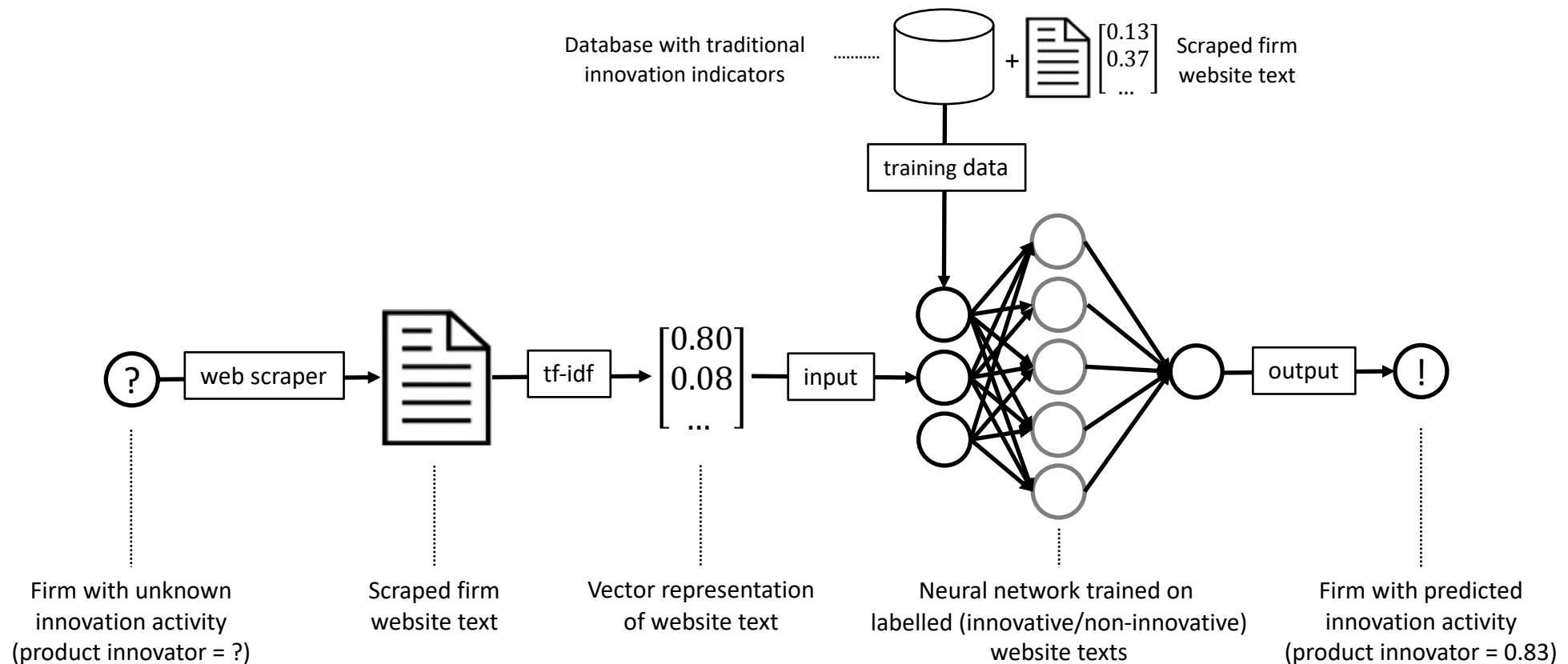
- ❖ Firms want us (as customers) to know how innovative they are.
- ❖ Almost all (significant) firms have websites nowadays.
- ❖ Important medium to tender and promote services and products.
- ❖ Firms have an incentive to keep their websites up-to-date.
- ❖ Firm website content is often related to product, service, and organisational innovations.

Using this openly accessible data to create timely and granular firm-level innovation indicators.

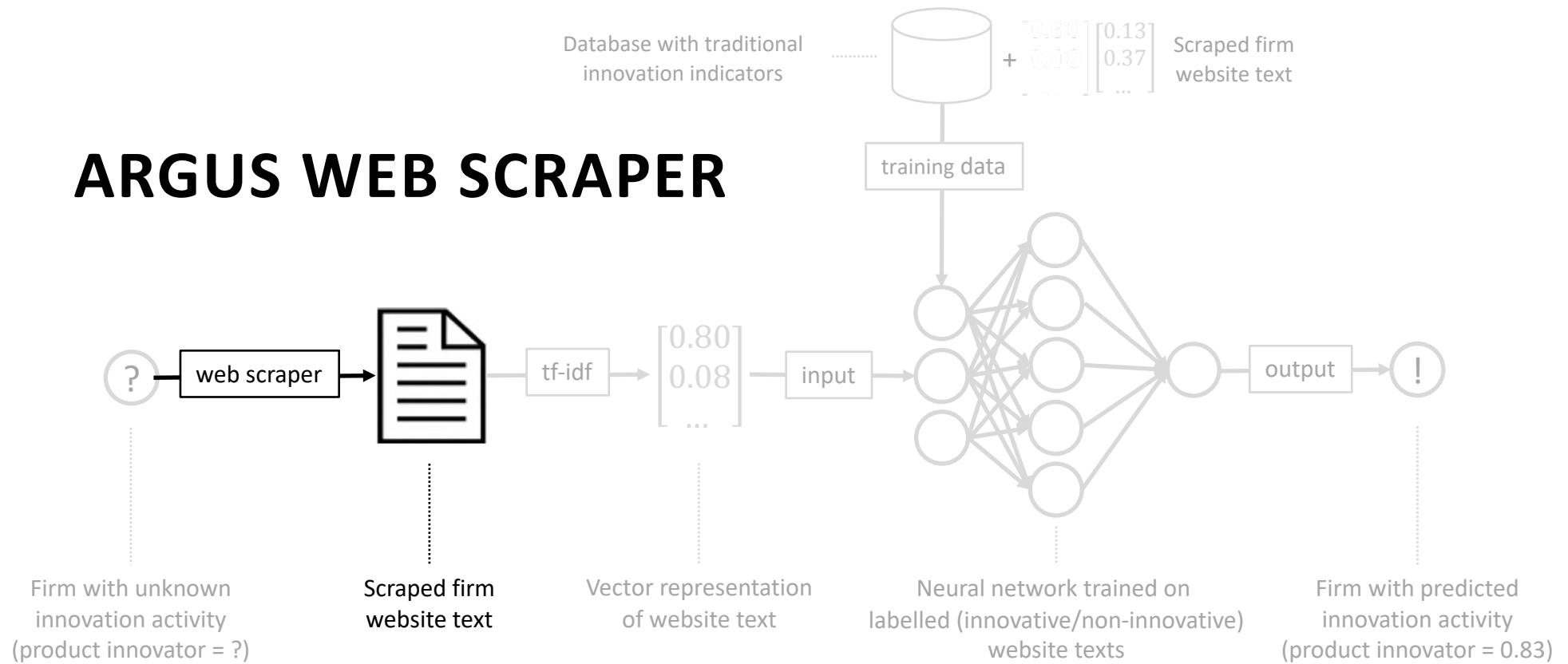
WEB MINING OF FIRM WEBSITES: ANALYSIS FRAMEWORK



DEEP LEARNING OF WEBSITE TEXT: FRAMEWORK

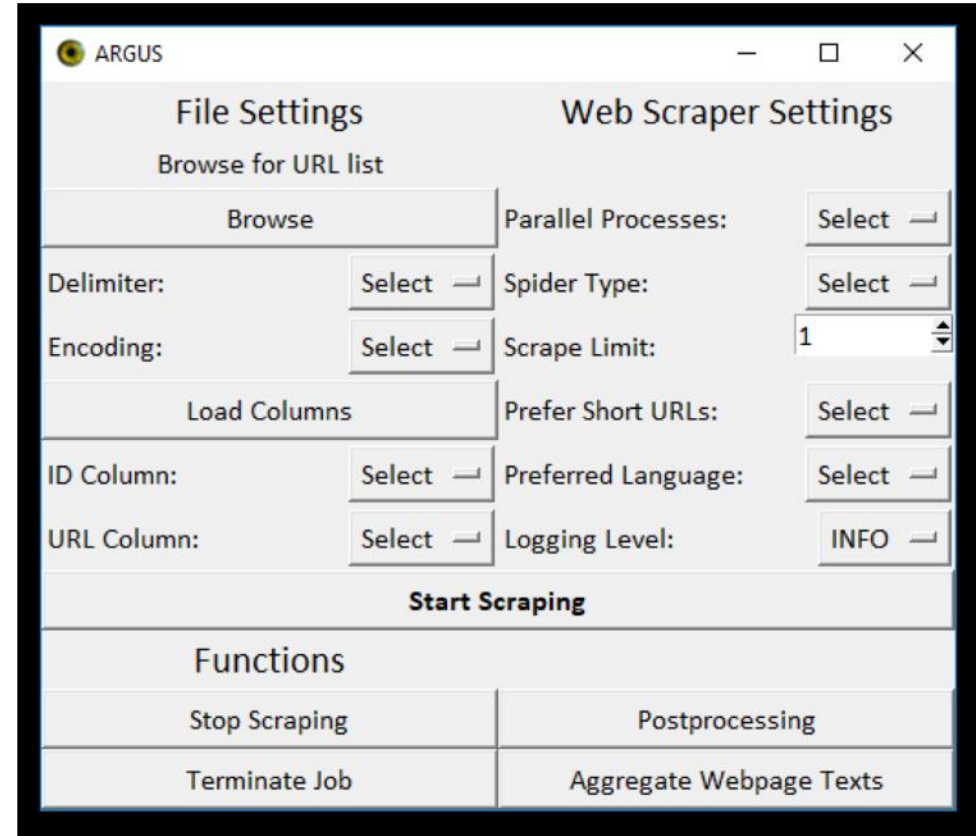


ARGUS WEB SCRAPER

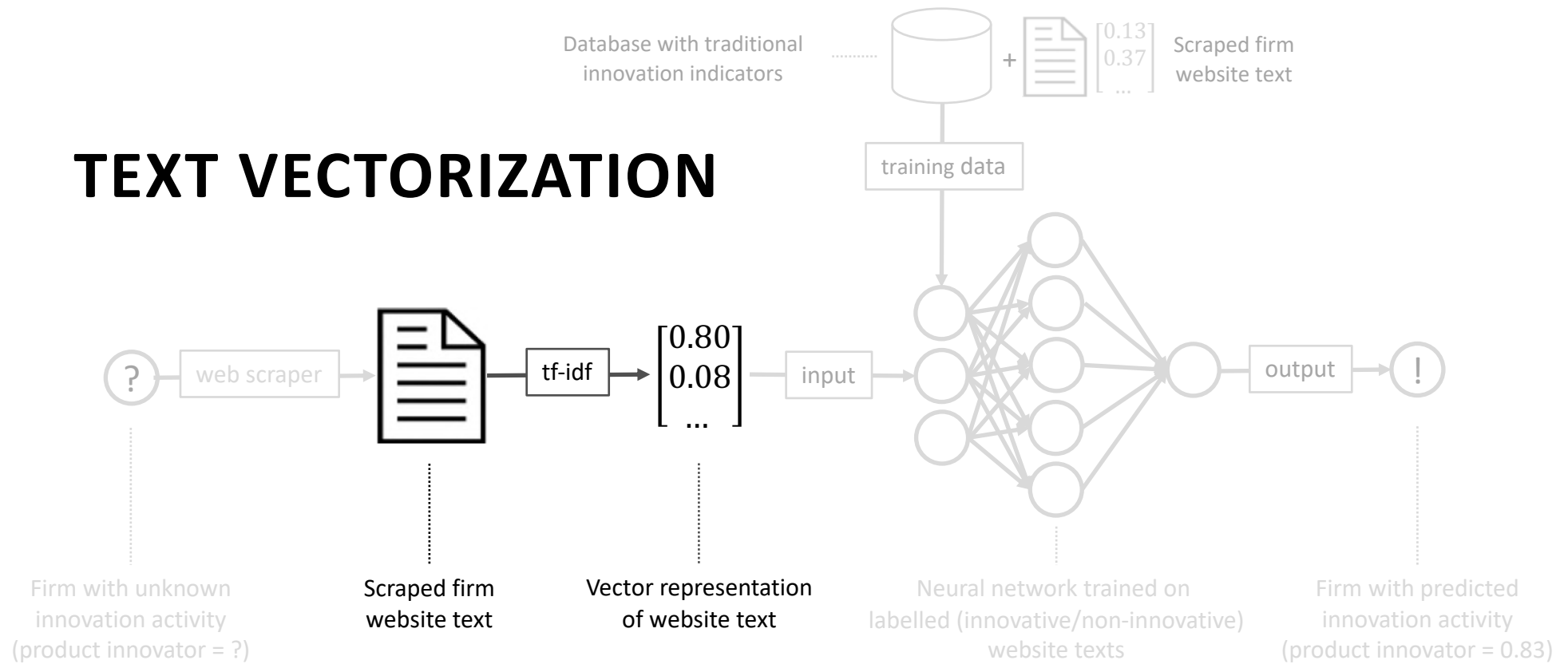


ARGUS WEB SCRAPER: OVERVIEW

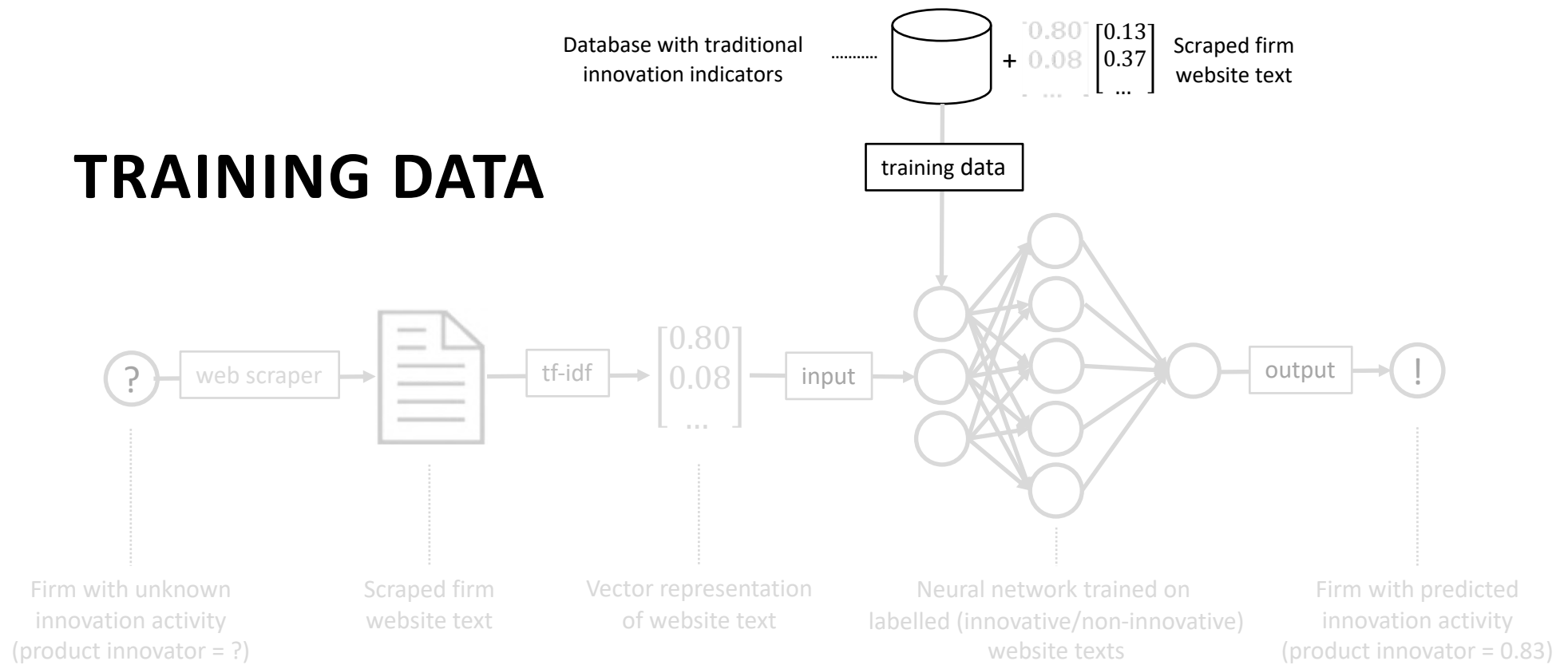
- ❖ Free and easy-to-use web scraping tool with GUI
[\[github.com/datawizard1337/ARGUS\]](https://github.com/datawizard1337/ARGUS)
- ❖ ~1,000 webpages per minute per CPU core
- ❖ ~5 days for 50m webpages on an office-grade PC
- ❖ Scraping of hyperlinks and texts



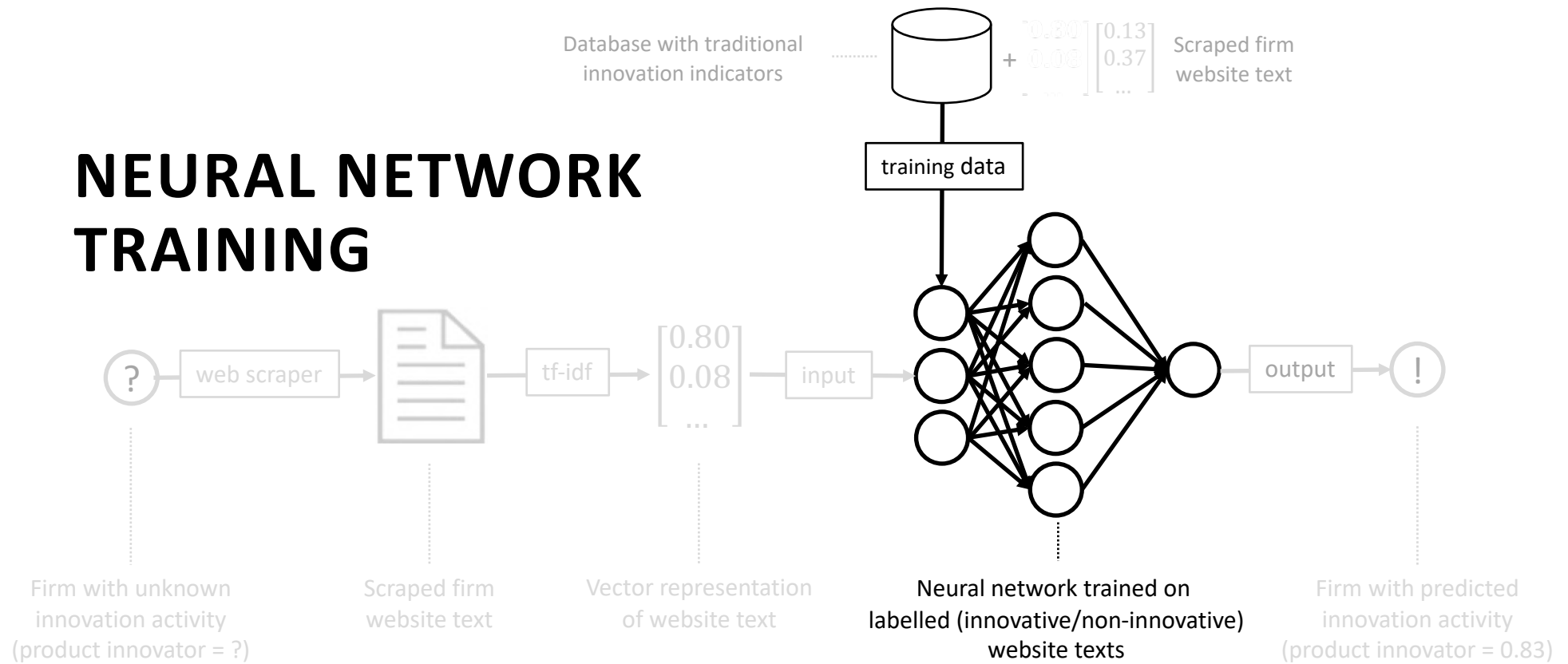
TEXT VECTORIZATION



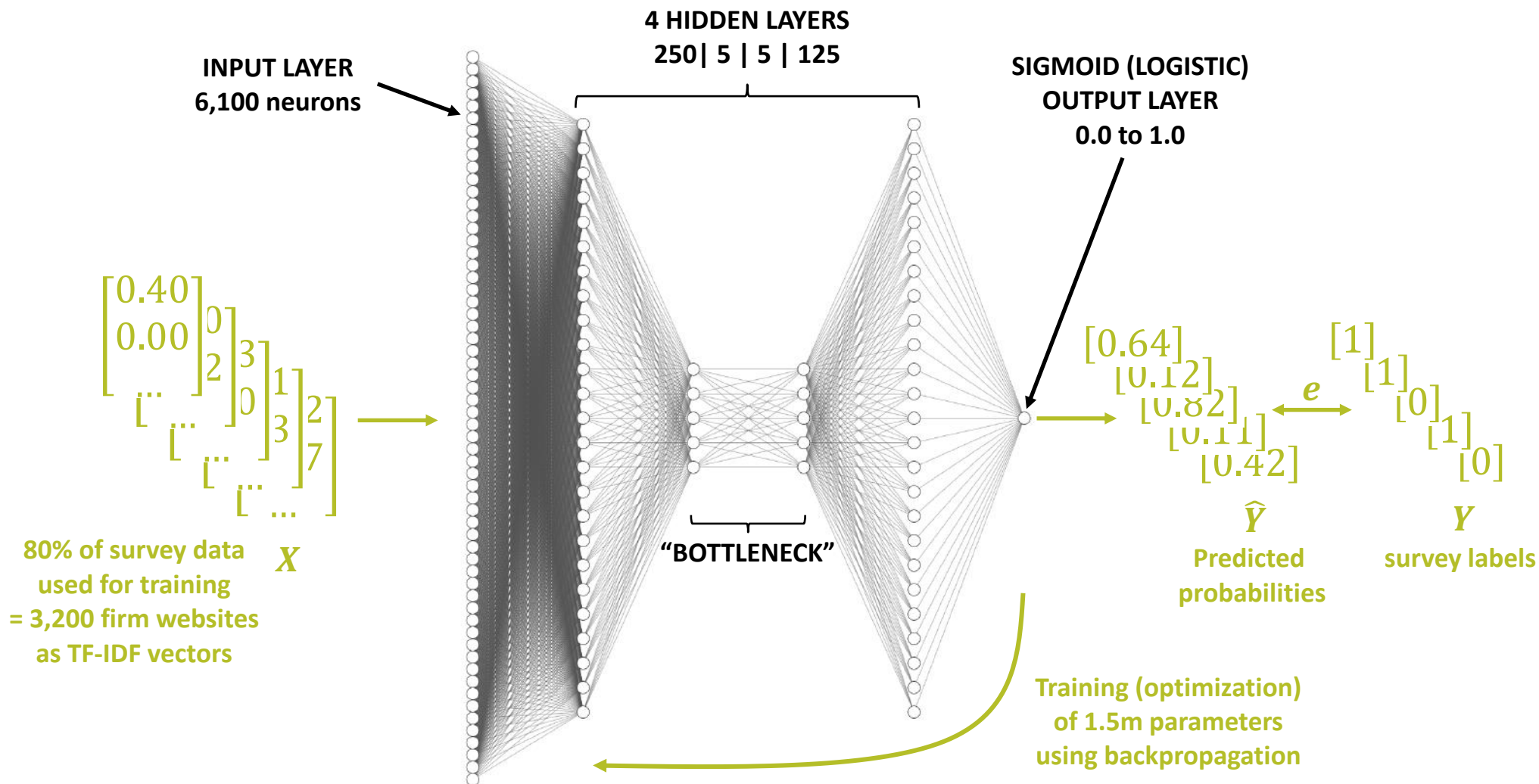
TRAINING DATA



NEURAL NETWORK TRAINING



NEURAL NETWORK TRAINING: UNDERCOMPLETE AUTOENCODER-LIKE ARCHITECTURE



NEURAL NETWORK TRAINING: PREDICTION PERFORMANCE

81% of firms predicted “non-innovative”
are actually “non-innovative”

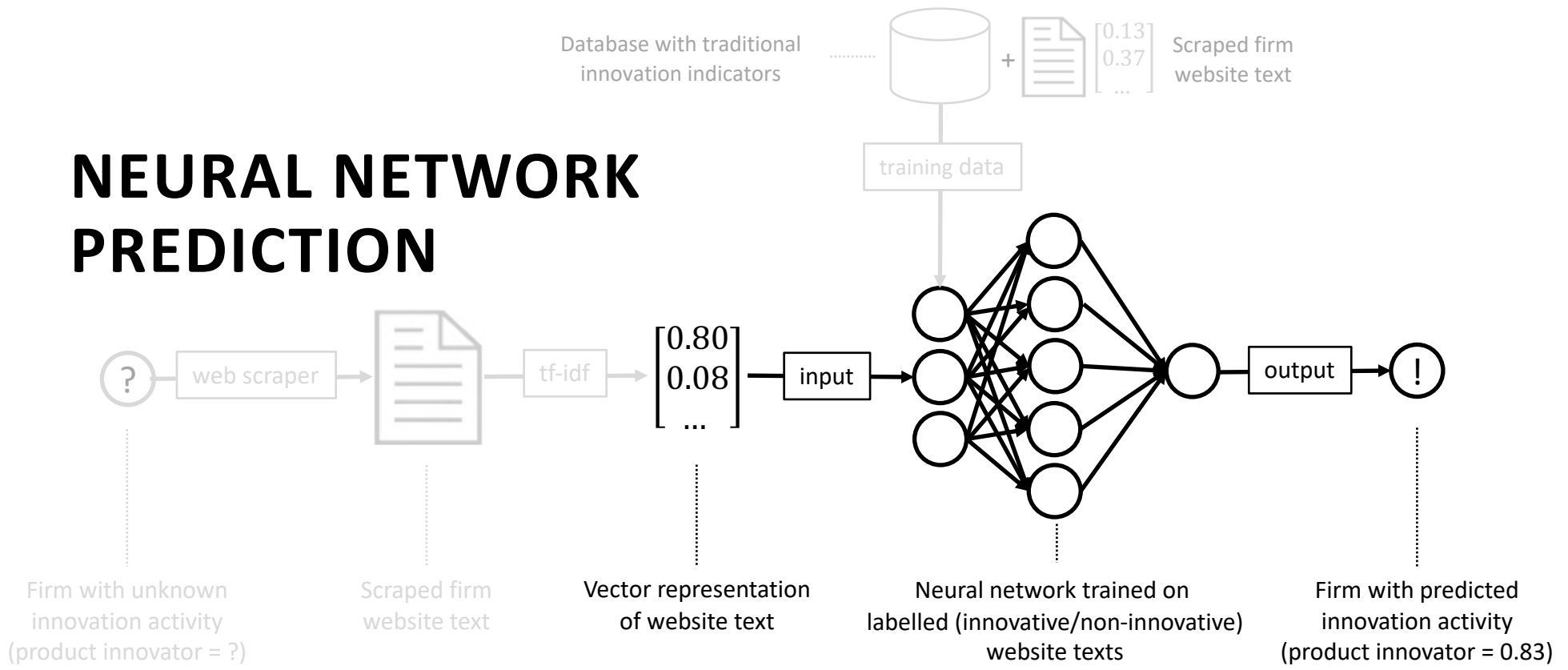
91% of actually “non-innovative”
firms are recovered

Composite measure
of precision and recall

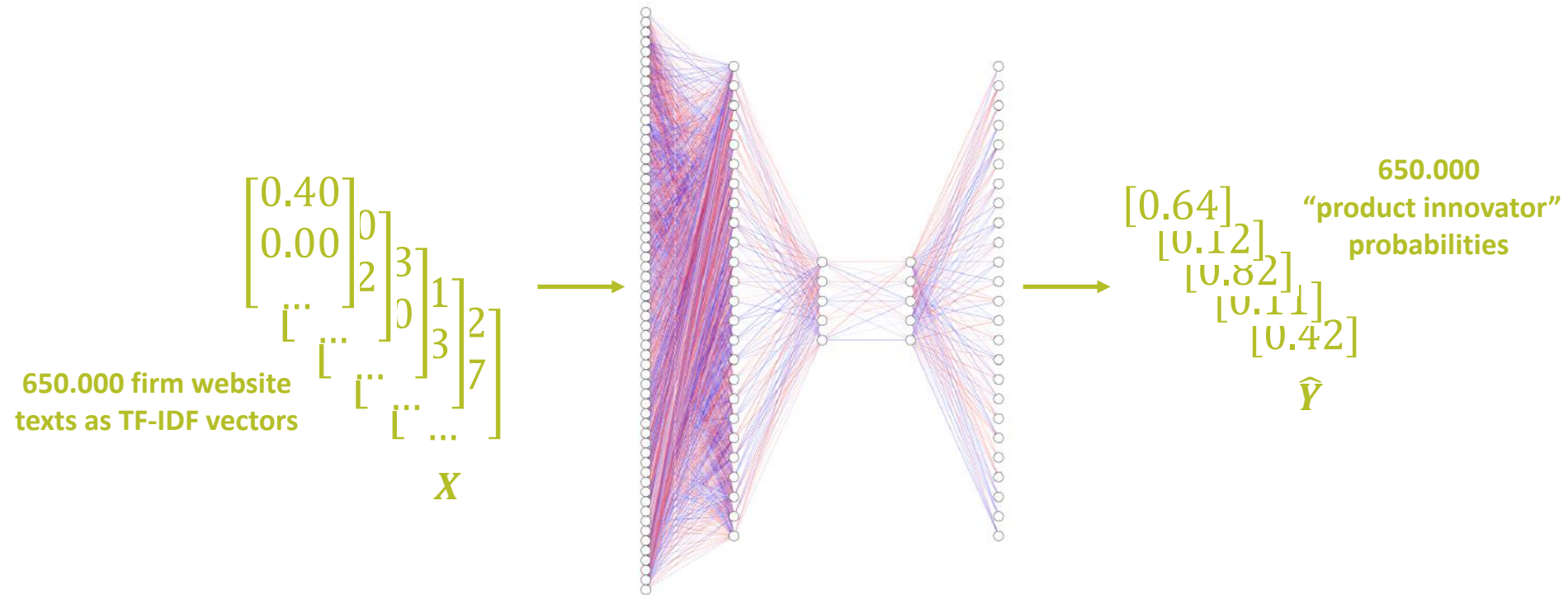
label	precision	recall	f1-score	support
non-innovative	0.81	0.91	0.86	429
innovative	0.81	0.64	0.71	255
avg / total	0.81	0.81	0.80	684

20% of survey data
as test set

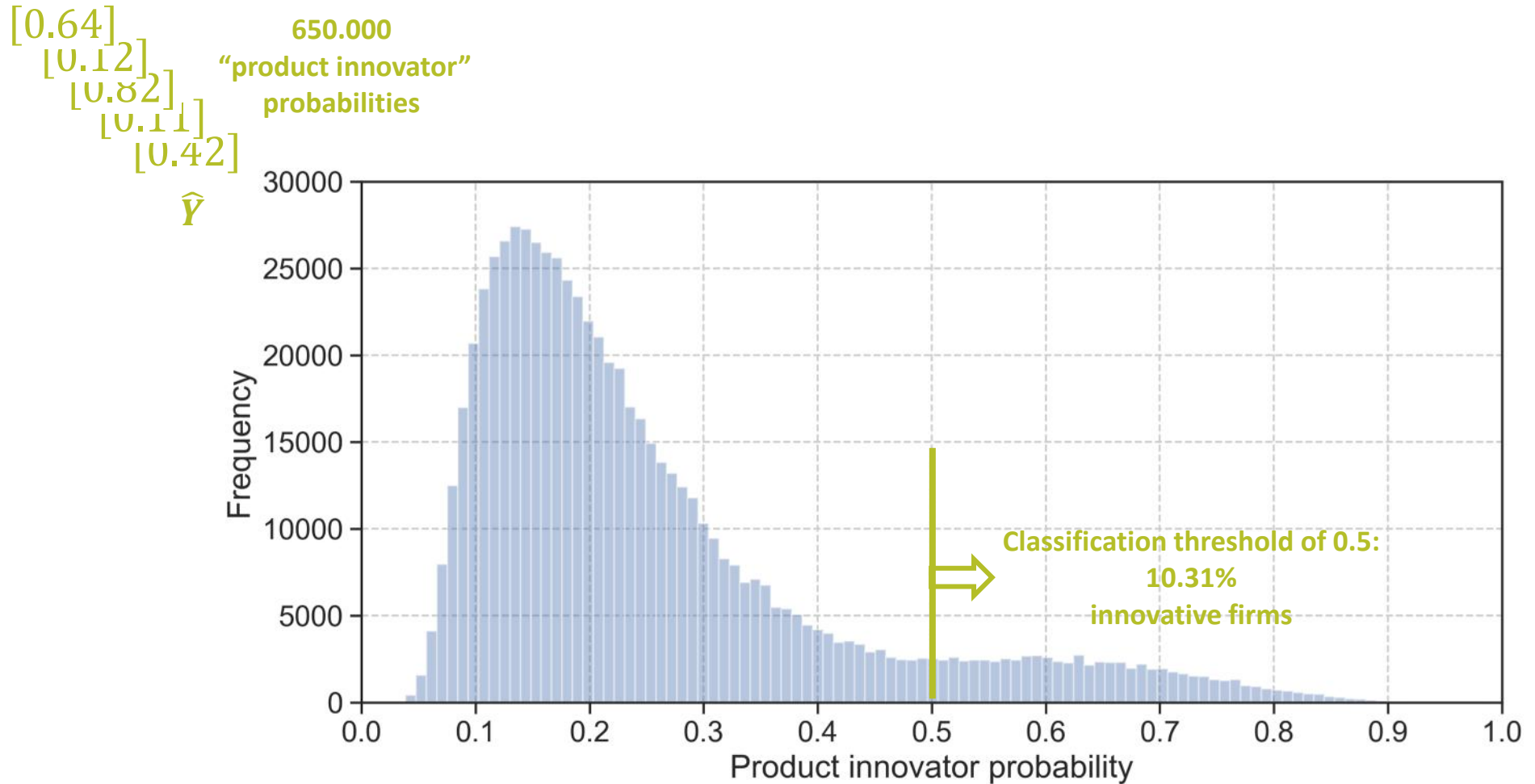
NEURAL NETWORK PREDICTION



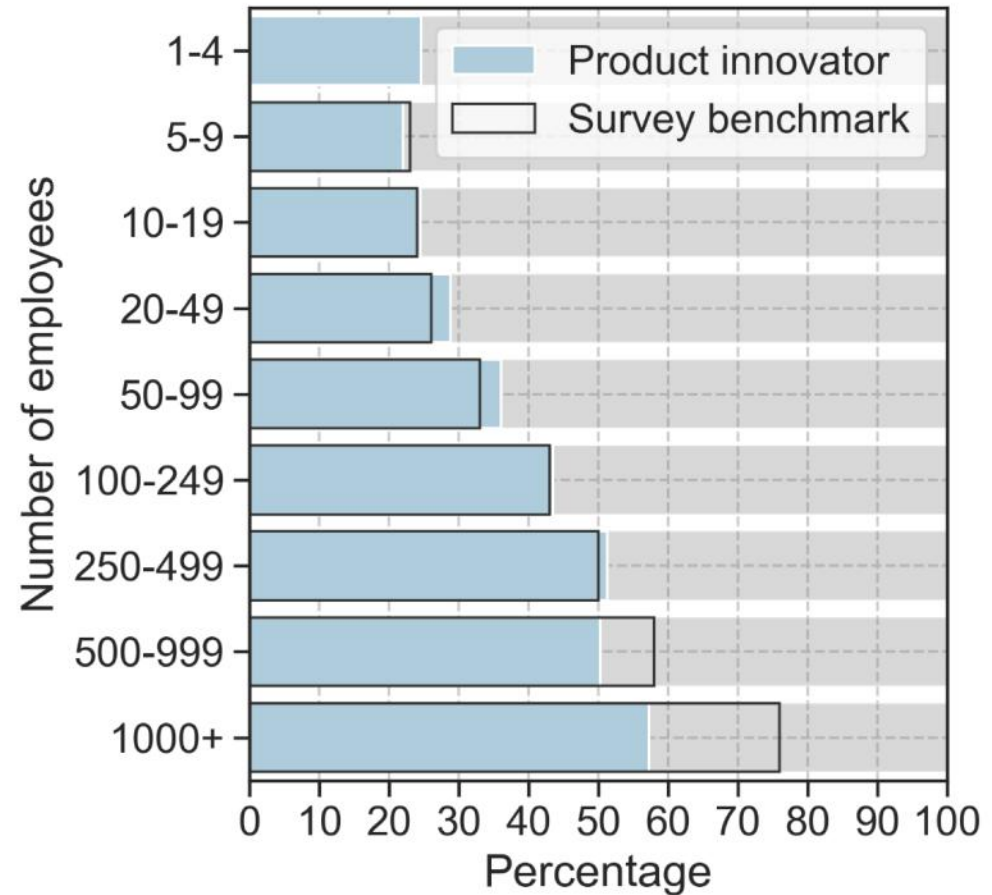
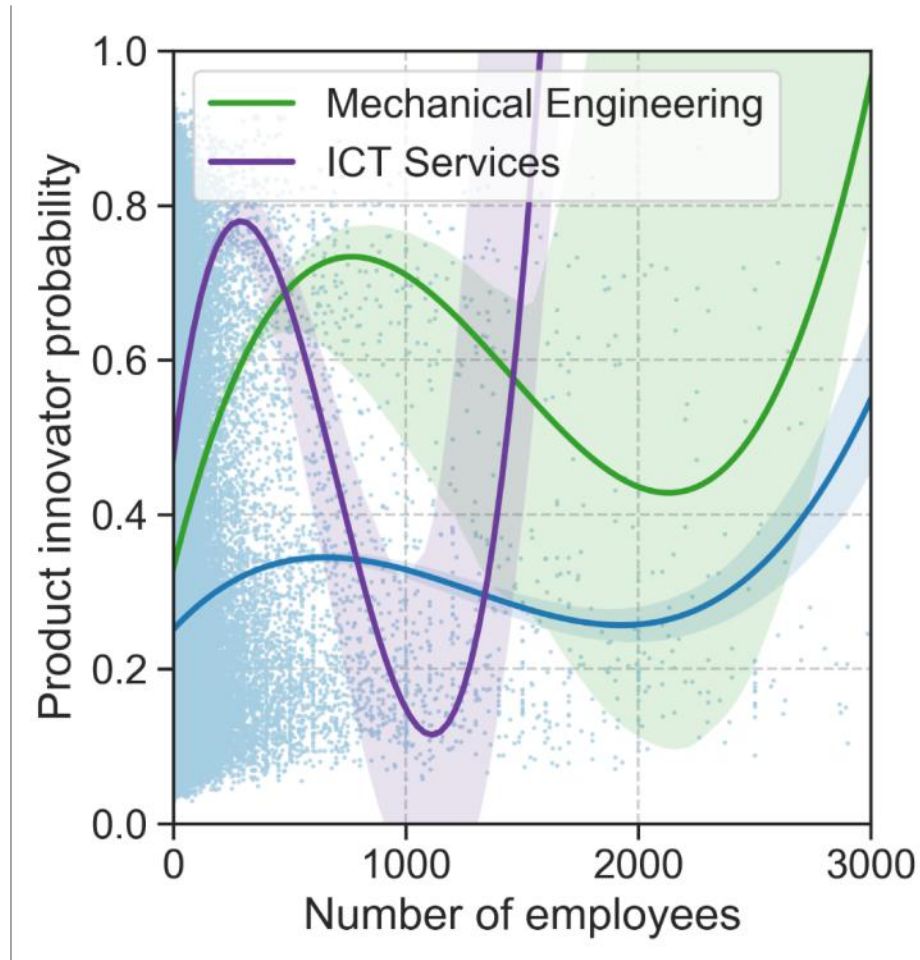
NEURAL NETWORK PREDICTIONS: OUT-OF-SAMPLE INNOVATION PREDICTION



NEURAL NETWORK PREDICTIONS: A CONTINUOUS INNOVATION INDICATOR

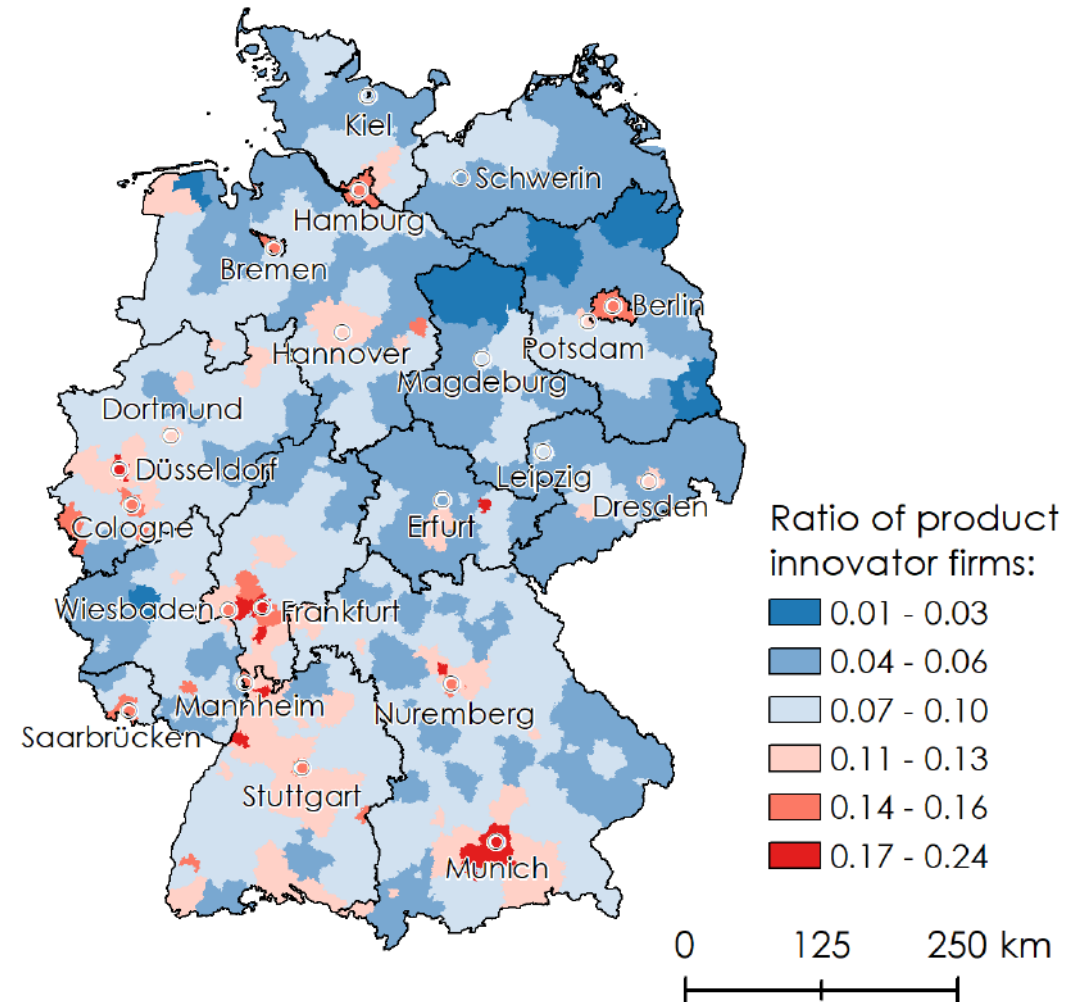


NEURAL NETWORK PREDICTIONS: INNOVATIVENESS BY FIRM SIZE

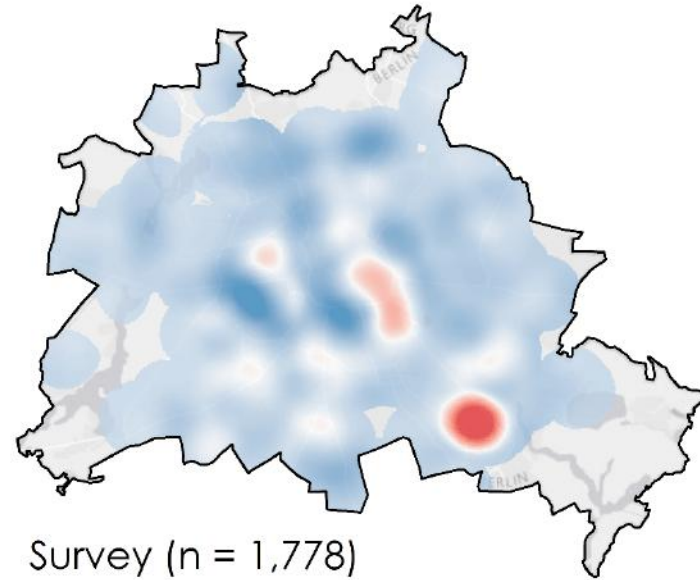
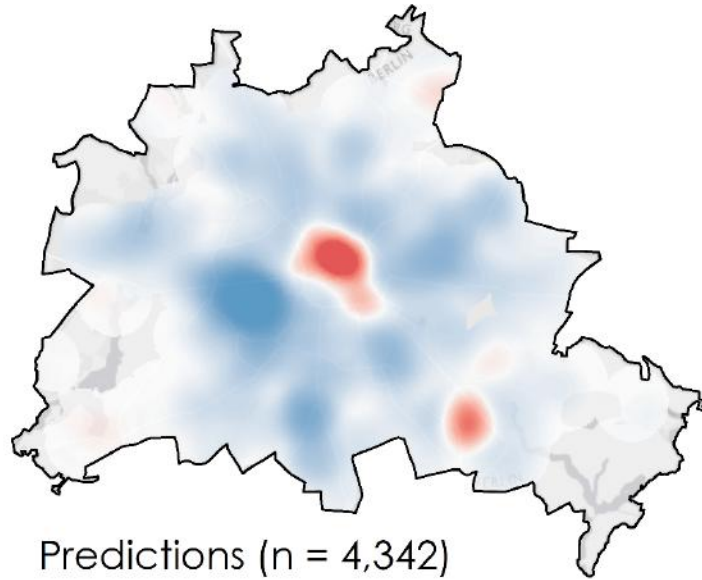


NEURAL NETWORK PREDICTIONS: INNOVATIVE DISTRICTS

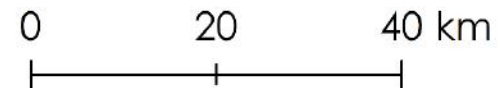
- ❖ Higher shares of product innovators in urban areas
 - Correlation with population density: 0.61
- ❖ Lower shares of product innovators in East and North



NEURAL NETWORK PREDICTIONS: MICROGEOGRAPHIC PATTERN



Dominant type of firm:



FUTURE DIRECTIONS: SUPPORTING BIG DATA BASED POLICY MAKING

- ❖ **Industry-specific prediction models**
- ❖ **Developing further web-based innovation indicators**
- ❖ **Formalization and dissemination of a coherent methodology**
- ❖ **Building up a panel database of web data**
- ❖ **Application in policy evaluation projects**
- ❖ **Investigating microgeographic diffusion of innovation and technology**

THANKS!



zew.de/en/team/jki



github.com/datawizard1337



twitter.com/jan_kinne