

PRACE SHAPE MOBILDEV – DATAMIN PROJESİ

Hazırlayan : Filiz Aksoy

- ▶ 18 yıllık deneyim
- ▶ 250 iş ortacı
- ▶ 10 servis sağlayıcı
- ▶ Digital pazarlama çözümleri
- ▶ Ürünler
 - ▶ İVT (İzinli veri tabanı),
 - ▶ Toplu SMS (Bulk SMS),
 - ▶ Toplu Email,
 - ▶ Anket Yönetimi vs
- ▶ Ar-ge : Datamin (TÜBİTAK- AGY100 - 03)

MOBİLDEV

▶ **Kişisel verilerin kategorize edilmesi**

➤ standart formatlı dokümanlar

* Kimlik, ehliyet, pasaport, askerlik belgesi, tapu, sicil kaydı....

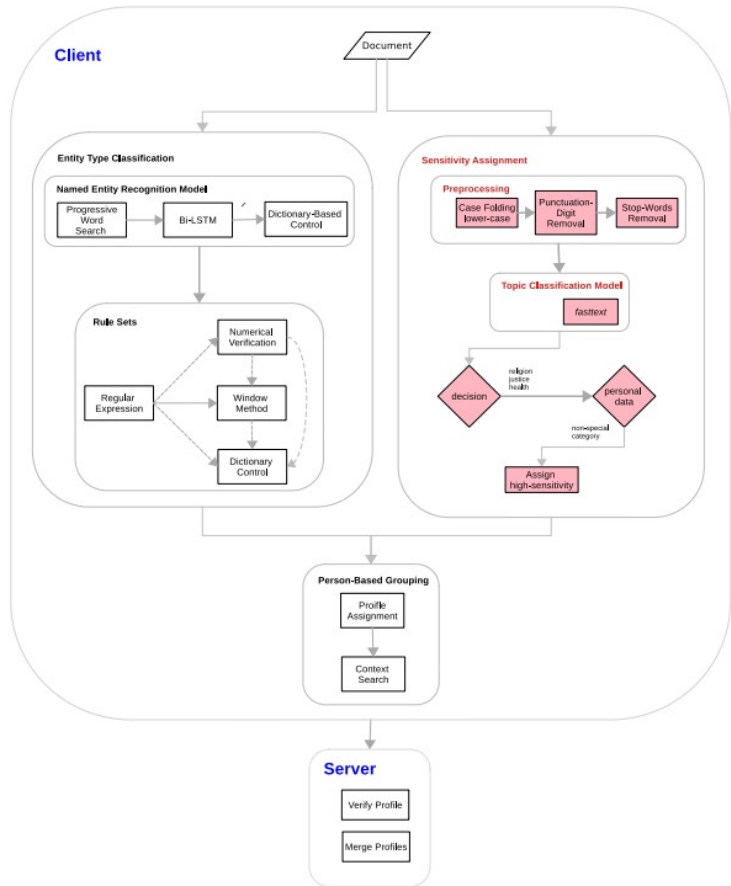
➤ Standart olmayan dokümanlar

* text, excel, pdf, word vs...

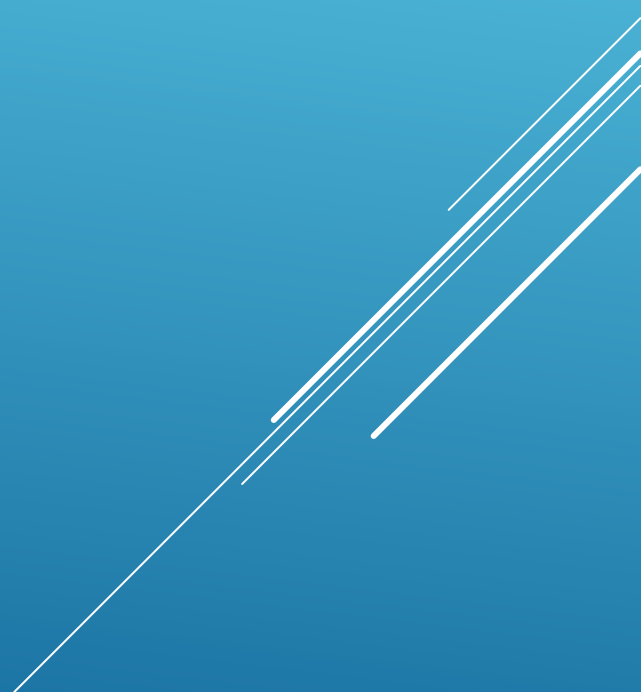
DATAMIN

- ▶ NER (Named Entity Recognition)
- ▶ LSTM (Long-Short Term Memory)
- ▶ Reg-ex (Regular expression)
- ▶ CNN
- ▶ Look-up table
- ▶ Dictionary
- ▶ Fasttext vektor uzayi & Bidirectional Encoder Representations Transformers (BERT)
- ▶ OCR - OpenCv

YONTEMLER



GENEL MIMARI YAPI



- ▶ Sistemdeki her düğüm, NVLink aracılığıyla bağlanan dört NVIDIA Tesla V100 Genel Amaçlı Grafik İşlem Birimi (GPGPU) içeriyor.
- ▶ Her GPU'nun 32 GB belleği ve her düğümün toplam 384 GB sistem belleği var.
- ▶ Düğümler arasındaki ara bağlantı ağı, Gelişmiş Veri Hızı InfiniBand aracılığıyla sağlanır. Donanım platformu ağırlıklı olarak fasttext ve BERT tabanlı modellerin eğitim aşaması için kullanılmıştır.

DATASET KATEGORI

Class	Number of instances
Politics	2549
Economy	3962
Culture	1852
Health	3078
Sport	10679
Technology	1471
Religion	4914
Justice	2574

MODEL TRAINING & MODEL TEST EXECUTION

		<i>fasttext</i>	BERT-based
Model Training	duration	1 minute 35 sec	7 hour 45 minutes 56 sec
	memory usage	1 GB	14 GB
	hardware	GPU Nvidia Tesla V4, 256 GB RAM	
	duration of an instance	1.6 ms	569 ms
Model Test Execution	hardware	CPU 8 Core, 16 GB RAM	
	model size	6.4 MB	438 MB
	memory usage	12 MB	987 MB
	CPU level	1	46

RESULT

Class	Model	Precision (%)	Recall (%)	F-1 Measure (%)
Technology	<i>fasttext</i>	91.88	99.32	95.45
	BERT-based	92.16	95.27	93.69
Culture	<i>fasttext</i>	94.87	100.00	97.36
	BERT-based	94.36	99.46	96.84
Justice	<i>fasttext</i>	87.85	87.85	87.85
	BERT-based	96.58	91.50	93.97
Economy	<i>fasttext</i>	94.69	98.74	96.67
	BERT-based	95.06	92.19	93.61
Health	<i>fasttext</i>	87.14	99.03	92.71
	BERT-based	83.01	96.75	89.36
Religion	<i>fasttext</i>	96.45	83.10	89.28
	BERT-based	97.99	89.41	93.50
Sport	<i>fasttext</i>	99.43	97.28	98.34
	BERT-based	99.81	98.69	99.25
Politics	<i>fasttext</i>	95.47	99.22	97.31
	BERT-based	91.82	96.86	94.27
Macro-avg ± st	<i>fasttext</i>	93.47 ± 1.50	95.57 ± 2.26	94.37 ± 1.40
	BERT-based	93.85 ± 1.82	95.02 ± 1.27	94.31 ± 1.00