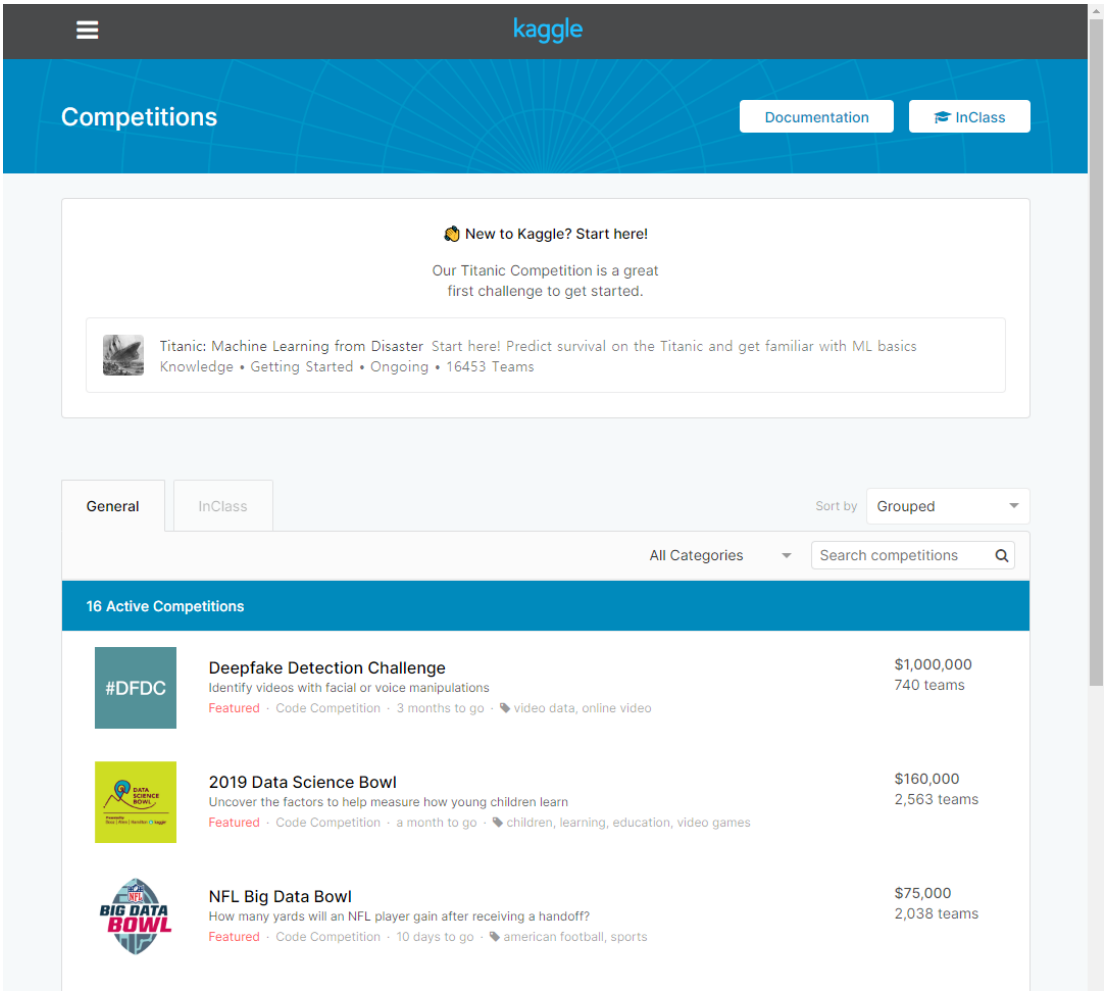


Chapter 04. 자연어처리 (Natural Language Processing)

관련 대회 소개

Kaggle



<https://www.kaggle.com/competitions>

Codalab


Codalab

Search CompetitionsMy CompetitionsHelpSign UpSign In

Competitions

Search...


Search



Evaluating grammatical error corrections
Organized by cnapoles
This "competition" contains different evaluation metrics commonly used for GEC and allows users to score their systems with these metrics ...

Nov 02, 2026- No end date


35 participants



ELEXIS Monolingual Word Sense Alignment Task
Organized by jmcrae
This is a competition to develop systems to predict alignment between senses of two monolingual dictionaries in 15 languages

Feb 03, 2020- No end date


10 participants



The Shared Task on Sarcasm Detection
Organized by debanjanhosh
Shared Task in Workshop on Figurative Language Processing

Jan 19, 2020-Mar 22, 2020


9 participants



The Second Shared Task on Metaphor Detection
Organized by cleong
Shared Task in Workshop on Figurative Language Processing

Jan 12, 2020-Mar 22, 2020

13 participants



NTIRE 2020 Image Deblurring Challenge - Track 2 on Smartphone
Organized by Radu
Our objectives are to gauge the s-o-t-a in example-based image deblurring on mobile devices, to promote research on this topic ...

Dec 27, 2019-Mar 15, 2020

3 participants

Join us on Github for contact & bug reports

About

Privacy and Terms

v1.5

<https://competitions.codalab.org/competitions/>

GLUE









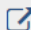
The General Language Understanding Evaluation (GLUE) benchmark is a collection of resources for training, evaluating, and analyzing natural language understanding systems. GLUE consists of:

- A benchmark of nine sentence- or sentence-pair language understanding tasks built on established existing datasets and selected to cover a diverse range of dataset sizes, text genres, and degrees of difficulty,
- A diagnostic dataset designed to evaluate and analyze model performance with respect to a wide range of linguistic phenomena found in natural language, and
- A public leaderboard for tracking performance on the benchmark and a dashboard for visualizing the performance of models on the diagnostic set.

The format of the GLUE benchmark is model-agnostic, so any system capable of processing sentence and sentence pairs and producing corresponding predictions is eligible to participate. The benchmark tasks are selected so as to favor models that share information across tasks using parameter sharing or other transfer learning techniques. The ultimate goal of GLUE is to drive research in the development of general and robust natural language understanding systems.

[PAPER](#) [STARTER CODE](#) [GROUP](#) [DIAGNOSTICS](#)

GLUE Tasks

Name	Download	More Info	Metric
The Corpus of Linguistic Acceptability			Matthew's Corr
The Stanford Sentiment Treebank			Accuracy
Microsoft Research Paraphrase Corpus			F1 / Accuracy
Semantic Textual Similarity Benchmark			Pearson-Spearman Corr
Quora Question Pairs			F1 / Accuracy
MultiNLI Matched			Accuracy
MultiNLI Mismatched			Accuracy
Question NLI			Accuracy
Recognizing Textual Entailment			Accuracy
Winograd NLI			Accuracy
Diagnostics Main			Matthew's Corr

DOWNLOAD DATA

<https://gluebenchmark.com/>

GLUE Leaderboard

Rank	Name	Model	URL	Score	CoLA	SST-2	MRPC	STS-B	QQP	MNLI-m	MNLI-mm	QNLI	RTE	
1	ERNIE Team - Baidu	ERNIE	Link	90.1	72.2	97.5	93.0/90.7	92.9/92.5	75.2/90.8	91.2	90.6	98.0	91.5	
2	Microsoft D365 AI & MSR AI & GATECHMT-DNN-SMART		Link	89.9	69.5	97.5	93.7/91.6	92.9/92.5	73.9/90.2	91.0	90.8	99.2	91.5	
3	T5 Team - Google	T5	Link	89.7	70.8	97.1	91.9/89.2	92.5/92.1	74.6/90.4	92.0	91.7	96.7	91.5	
+	4	王玮	ALICE v2 large ensemble (Alibaba DAMO NLP)	Link	89.5	71.3	97.1	93.9/91.9	93.0/92.5	74.8/91.0	90.7	90.4	99.2	91.5
5	XLNet Team	XLNet (ensemble)	Link	89.5	70.2	97.1	92.9/90.5	93.0/92.6	74.7/90.4	90.9	90.9	99.0	91.5	
6	ALBERT-Team Google Language	ALBERT (Ensemble)	Link	89.4	69.1	97.1	93.4/91.2	92.5/92.0	74.2/90.5	91.3	91.0	99.2	91.5	
7	Microsoft D365 AI & UMD	FreeLB-RoBERTa (ensemble)	Link	88.8	68.0	96.8	93.1/90.8	92.4/92.2	74.8/90.3	91.1	90.7	98.8	91.5	
8	Facebook AI	RoBERTa	Link	88.5	67.8	96.7	92.3/89.8	92.2/91.9	74.3/90.2	90.8	90.2	98.9	91.5	
9	Junjie Yang	HIRE-RoBERTa	Link	88.3	68.6	97.1	93.0/90.7	92.4/92.0	74.3/90.2	90.7	90.4	95.5	91.5	
+	10	Microsoft D365 AI & MSR AI	MT-DNN-ensemble	Link	87.6	68.4	96.5	92.7/90.3	91.1/90.7	73.7/89.9	87.9	87.4	96.0	91.5
11	GLUE Human Baselines	GLUE Human Baselines	Link	87.1	66.4	97.8	86.3/80.8	92.7/92.6	59.5/80.4	92.0	92.8	91.2	91.5	
12	Stanford Hazy Research	Snorkel MeTaL	Link	83.2	63.8	96.2	91.5/88.5	90.1/89.7	73.1/89.9	87.6	87.2	93.9	91.5	
13	XLM Systems	XLM (English only)	Link	83.1	62.9	95.6	90.7/87.1	88.8/88.2	73.2/89.8	89.1	88.5	94.0	91.5	
14	Zhuosheng Zhang	SemBERT	Link	82.9	62.3	94.6	91.2/88.3	87.8/86.7	72.8/89.8	87.6	86.3	94.6	91.5	
15	Danqi Chen	SpanBERT (single-task training)	Link	82.8	64.3	94.8	90.9/87.9	89.9/89.1	71.9/89.5	88.1	87.7	94.3	91.5	
16	Kevin Clark	BERT + BAM	Link	82.3	61.5	95.2	91.3/88.3	88.6/87.9	72.5/89.7	86.6	85.8	93.1	91.5	

GLUE의 리더보드에서는 다양한 최신 연구를 실시간으로 확인할 수 있다.

SuperGLUE











In the last year, new models and methods for pretraining and transfer learning have driven striking performance improvements across a range of language understanding tasks. The GLUE benchmark, introduced one year ago, offered a single-number metric that summarizes progress on a diverse set of such tasks, but performance on the benchmark has recently come close to the level of non-expert humans, suggesting limited headroom for further research.





















We take into account the lessons learnt from original GLUE benchmark and present SuperGLUE, a new benchmark styled after GLUE with a new set of more difficult language understanding tasks, improved resources, and a new public leaderboard.

PAPER

STARTER CODE

GROUP

DIAGNOSTICS







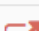

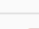
SuperGLUE Tasks				
Name	Identifier	Download	More Info	Metric
Broadcoverage Diagnostics	AX-b			Matthew's Corr
CommitmentBank	CB			Avg. F1 / Accuracy
Choice of Plausible Alternatives	COPA			Accuracy
Multi-Sentence Reading Comprehension	MultiRC			F1a / EM
Recognizing Textual Entailment	RTE			Accuracy
Words in Context	WIC			Accuracy
The Winograd Schema Challenge	WSC			Accuracy
BoolQ	BoolQ			Accuracy
Reading Comprehension with Commonsense Reasoning	ReCoRD			F1 / Accuracy
Winogender Schema Diagnostics	AX-g			Gender Parity / Accuracy

DOWNLOAD ALL DATA

<https://super.gluebenchmark.com/>

SuperGLUE Leaderboard

Leaderboard Version: 2.0

Rank	Name	Model	URL	Score	BoolQ	CB	COPA	MultiRC	ReCoRD	RTE	WiC	WSC	AX-b	AX-g
1	SuperGLUE Human Baselines	SuperGLUE Human Baselines		89.8	89.0	95.8/98.9	100.0	81.8/51.9	91.7/91.3	93.6	80.0	100.0	76.6	99.3/99.7
2	T5 Team - Google	T5		88.9	91.0	93.0/96.4	94.8	88.2/62.3	93.3/92.5	92.5	76.1	93.8	65.6	92.7/91.9
3	Facebook AI	RoBERTa		84.6	87.1	90.5/95.2	90.6	84.4/52.5	90.6/90.0	88.2	69.9	89.0	57.9	91.0/78.1
4	IBM Research AI	BERT-ml		73.5	84.8	89.6/94.0	73.8	73.2/30.5	74.6/74.0	84.1	66.2	61.0	29.6	97.8/57.3
5	SuperGLUE Baselines	BERT++		71.5	79.0	84.8/90.4	73.8	70.0/24.1	72.0/71.3	79.0	69.6	64.4	38.0	99.4/51.4
		BERT		69.0	77.4	75.7/83.6	70.6	70.0/24.1	72.0/71.3	71.7	69.6	64.4	23.0	97.8/51.7
		Most Frequent Class		47.1	62.3	21.7/48.4	50.0	61.1/0.3	33.4/32.5	50.3	50.0	65.1	0.0	100.0/50.0
		CBoW		44.5	62.2	49.0/71.2	51.6	0.0/0.5	14.0/13.6	49.7	53.1	65.1	-0.4	100.0/50.0
		Outside Best		-	80.4	-	84.4	70.4/24.5	74.8/73.0	82.7	-	-	-	-
-	Stanford Hazy Research	Snorkel [SuperGLUE v1.9]		-	-	88.6/93.2	76.2	76.4/36.3	-	78.9	72.1	72.6	47.6	-

SuperGLUE의 경우 아직 GLUE보다 많은 연구가 수행되지 않았다.

발전하는 과정을 관찰할 수도 있고, 직접 뛰어들어 좋은 성능을 뽑내볼 수도 있다.