



Universität Hamburg  
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Alexander Panchenko

---

**FROM UNSUPERVISED INDUCTION OF  
LINGUISTIC STRUCTURES FROM TEXT  
TOWARDS APPLICATIONS IN DEEP  
LEARNING**

## Shared task on word sense induction

# A shared task on WSI

- An **ACL SIGSLAV** sponsored shared task on **word sense induction (WSI)** for the Russian language.
- **More details:** <http://russe.nlpub.org/2018/wsi>



# A lexical sample WSI task

- **Target word**, e.g. “bank”.

# A lexical sample WSI task

- **Target word**, e.g. “bank”.
- **Contexts** where the word occurs, e.g.:
  - “river **bank** is a slope beside a body of water”
  - “**bank** is a financial institution that accepts deposits”
  - “Oh, the **bank** was robbed. They took about a million dollars.”
  - “**bank** of Elbe is a good and popular hangout spot complete with good food and fun”

# A lexical sample WSI task

- **Target word**, e.g. “bank”.
- **Contexts** where the word occurs, e.g.:
  - “river **bank** is a slope beside a body of water”
  - “**bank** is a financial institution that accepts deposits”
  - “Oh, the **bank** was robbed. They took about a million dollars.”
  - “**bank** of Elbe is a good and popular hangout spot complete with good food and fun”
- You need to **group the contexts by senses**:
  - “river **bank** is a slope beside a body of water”
  - “**bank** of Elbe is a good and popular hangout spot complete with good food and fun”
  - “**bank** is a financial institution that accepts deposits”
  - “Oh, the **bank** was robbed. They took about a million dollars.”

# Dataset based on Wikipedia

Dataset	Type	Inventory	Corpus	Split	# of words	# of senses	Avg. # of senses	# of contexts
wiki-wiki	main	Wikipedia	Wikipedia	train	4	8	2.0	439
wiki-wiki	main	Wikipedia	Wikipedia	test	5	12	2.4	539
bts-rnc	main	BTS	RNC	train	30	96	3.2	3491
bts-rnc	main	BTS	RNC	test	51	153	3.0	6556
active-dict	main	Active Dict.	Active Dict.	train	85	312	3.7	2073
active-dict	main	Active Dict.	Active Dict.	test	168	555	3.3	3729
active-rnc	additional	Active Dict.	RNC	train	20	71	3.6	1829
active-rutenten	additional	Active Dict.	ruTenTen <sup>12</sup>	train	21	71	3.4	3671
bts-rutenten	additional	BTS	ruTenTen	train	11	25	2.3	956

# Dataset based on RNC

Dataset	Type	Inventory	Corpus	Split	# of words	# of senses	Avg. # of senses	# of contexts
wiki-wiki	main	Wikipedia	Wikipedia	train	4	8	2.0	439
wiki-wiki	main	Wikipedia	Wikipedia	test	5	12	2.4	539
bts-rnc	main	BTS	RNC	train	30	96	3.2	3491
bts-rnc	main	BTS	RNC	test	51	153	3.0	6556
active-dict	main	Active Dict.	Active Dict.	train	85	312	3.7	2073
active-dict	main	Active Dict.	Active Dict.	test	168	555	3.3	3729
active-rnc	additional	Active Dict.	RNC	train	20	71	3.6	1829
active-rutenten	additional	Active Dict.	ruTenTen <sup>12</sup>	train	21	71	3.4	3671
bts-rutenten	additional	BTS	ruTenTen	train	11	25	2.3	956



# Dataset based on dictionary glosses

Dataset	Type	Inventory	Corpus	Split	# of words	# of senses	Avg. # of senses	# of contexts
wiki-wiki	main	Wikipedia	Wikipedia	train	4	8	2.0	439
wiki-wiki	main	Wikipedia	Wikipedia	test	5	12	2.4	539
bts-rnc	main	BTS	RNC	train	30	96	3.2	3491
bts-rnc	main	BTS	RNC	test	51	153	3.0	6556
active-dict	main	Active Dict.	Active Dict.	train	85	312	3.7	2073
active-dict	main	Active Dict.	Active Dict.	test	168	555	3.3	3729
active-rnc	additional	Active Dict.	RNC	train	20	71	3.6	1829
active-rutenten	additional	Active Dict.	ruTenTen <sup>12</sup>	train	21	71	3.4	3671
bts-rutenten	additional	BTS	ruTenTen	train	11	25	2.3	956

# A sample from the *wiki-wiki* dataset

word	articles	sense
белка	кавказская белка; обыкновенная белка; японская белка; капская земляная белка; аризонская белка; ...	рыжая, шустрая, дерево, вскарабкаться, прыгнуть
белка	домен белка; биосинтез белка; фолдинг белка; институт белка ран; сигнальная функция белка; ...	желток, пища, углевод, рацион, жир
белка	белка и стрелка; белка и стрелка (мюзикл); белка и стрелка. лунные приключения; ...	космос, полет, животные, первые, советские

# A sample from the *wiki-wiki* dataset

Rank	Team	ARI (public)	ARI (private)
1	jamsic	1.0000 (1)	0.9625 (1)
2	akutuzov (Kutuzov, 2018)	0.9823 (2)	0.7096 (2)
3	ezhick179 (Arefyev et al., 2018)	1.0000 (1)	0.6586 (3)
-	<i>akapustin</i>	0.6520 (6)	0.6459 (4)
-	<i>aby2s</i>	1.0000 (1)	0.5889 (5)
-	<i>bokan</i>	0.7587 (5)	0.5530 (6)
*	<b>AdaGram (Bartunov et al, 2016)</b>	<b>0.6278 (7)</b>	<b>0.5275 (7)</b>
4	Pavel (Arefyev et al., 2018)	0.9649 (3)	0.4827 (8)
5	eugenys	0.0115 (12)	0.4377 (9)
6	mikhal	1.0000 (1)	0.4109 (10)
7	fogside	0.6520 (6)	0.3958 (11)

# A sample from the *wiki-wiki* dataset

Rank	Team	ARI (public)	ARI (private)
1	<b>embeddings + sense induction</b>		0.9625 (1)
2	<b>embeddings + affinity propagation</b>		0.7096 (2)
3	<b>embeddings + agglomerative clustering</b>		0.6586 (3)
-	<i>akapustin</i>	0.6520 (6)	0.6459 (4)
-	<i>aby2s</i>	1.0000 (1)	0.5889 (5)
-	<i>bokan</i>	0.7587 (5)	0.5530 (6)
*	<b>AdaGram (Bartunov et al, 2016)</b>	<b>0.6278 (7)</b>	<b>0.5275 (7)</b>
4	Pavel (Arefyev et al., 2018)	0.9649 (3)	0.4827 (8)
5	eugenys	0.0115 (12)	0.4377 (9)
6	mikhal	1.0000 (1)	0.4109 (10)
7	fogside	0.6520 (6)	0.3958 (11)