Splitting and Identifying Czech Compounds: A Pilot Study

Emil Svoboda

Charles University, Faculty of Mathematics and Physics



Joint work with Magda Ševčíková for DeriMo 2021

- Compounding and its challenges in Czech
- Data set and evaluation description
- Our solutions of identifying and splitting compounds
 - Baseline
 - o Interlexical Matrices of Likeness based heuristic method
 - Czech Compound Splitter
- Results and conclusion

 $v \check{z} dy + z elen \check{y} \rightarrow v \check{z} dy z elen \check{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words"

- Czech, where compounding is underresearched computationally
- graphical words only:
 - $\circ~$ strings of graphical symbols unbroken by whitespace or special characters

 $v\ddot{z}dy$ + $zelen\acute{y}$ ightarrow $v\ddot{z}dyzelen\acute{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words"

This paper focuses on:

- Czech, where compounding is underresearched computationally
- graphical words only:
 - $\circ~$ strings of graphical symbols unbroken by whitespace or special characters

 $v\check{z}dy$ + $zelen\check{y}$ ightarrow $v\check{z}dyzelen\check{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words" This paper focuses on:

- Czech, where compounding is underresearched computationally
- graphical words only:
 - $\circ~$ strings of graphical symbols unbroken by whitespace or special characters

 $v\check{z}dy + zelen\check{y}
ightarrow v\check{z}dyzelen\check{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words"

This paper focuses on:

- Czech, where compounding is underresearched computationally
- graphical words only:
 - $\circ~$ strings of graphical symbols unbroken by whitespace or special characters

 $v\check{z}dy + zelen\check{y} \rightarrow v\check{z}dyzelen\check{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words" This paper focuses on:

- Czech, where compounding is underresearched computationally
- graphical words only:
 - $\circ\,$ strings of graphical symbols unbroken by whitespace or special characters

 $v\check{z}dy + zelen\check{y}
ightarrow v\check{z}dyzelen\check{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words" This paper focuses on:

- Czech, where compounding is underresearched computationally
- graphical words only:
 - $\circ\,$ strings of graphical symbols unbroken by whitespace or special characters

 $v\check{z}dy + zelen\check{y} \rightarrow v\check{z}dyzelen\check{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words" This paper focuses on:

- Czech, where compounding is underresearched computationally
- graphical words only:

 $\circ\,$ strings of graphical symbols unbroken by whitespace or special characters

 $v\check{z}dy$ + $zelen\check{y}$ ightarrow $v\check{z}dyzelen\check{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words" This paper focuses on:

- Czech, where compounding is underresearched computationally
- graphical words only:
 - $\circ~$ strings of graphical symbols unbroken by whitespace or special characters

 $v\check{z}dy$ + $zelen\check{y}$ ightarrow $v\check{z}dyzelen\check{y}$ always.ADV green.ADJ evergreen.ADJ

• the two or more input lexemes are "parents" or "parent words" This paper focuses on:

- Czech, where compounding is underresearched computationally
- graphical words only:
 - $\circ~$ strings of graphical symbols unbroken by whitespace or special characters

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - \circ if **W** is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - $\circ\,$ return linguistically valid parent lemmas of W, separated by spaces.

Compound identification and compound splitting

Compound identification:

• Binary classification

- Input graphical word **W** (assumed to be a Czech lexeme):
 - \circ if **W** is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - $\circ~$ return linguistically valid parent lemmas of W, separated by spaces.

Compound identification and compound splitting

Compound identification:

• Binary classification

- Input graphical word **W** (assumed to be a Czech lexeme):
 - \circ if **W** is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - $\circ~$ return linguistically valid parent lemmas of W, separated by spaces.

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - $\circ\,$ if ${\bm W}$ is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - $\circ\,$ return linguistically valid parent lemmas of W, separated by spaces.

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - $\circ~$ if \boldsymbol{W} is a Czech compound, return 1;
 - o else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - $\circ\,$ return linguistically valid parent lemmas of W_{i} separated by spaces.

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - $\circ~$ if \boldsymbol{W} is a Czech compound, return 1;
 - \circ else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - $\circ\,$ return linguistically valid parent lemmas of W, separated by spaces.

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - $\circ~$ if \boldsymbol{W} is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - $\circ~$ return linguistically valid parent lemmas of $\boldsymbol{W},$ separated by spaces.

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - if **W** is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - \circ return *linguistically valid* parent *lemmas* of W, separated by spaces.

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - if **W** is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word ${\bf W}$ (assumed to be a Czech compound):

o return linguistically valid parent lemmas of W, separated by spaces.

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - if **W** is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word ${\bf W}$ (assumed to be a Czech compound):
 - $\circ~$ return linguistically valid parent lemmas of \boldsymbol{W} , separated by spaces.

Compound identification and compound splitting

Compound identification:

- Binary classification
- Input graphical word **W** (assumed to be a Czech lexeme):
 - if **W** is a Czech compound, return 1;
 - else, return 0.

Compound splitting:

- An open-ended task
- Input graphical word **W** (assumed to be a Czech compound):
 - return linguistically valid parent lemmas of W, separated by spaces.

What makes the tasks difficult?

What makes the tasks difficult?

What makes the tasks difficult?

 $krv\text{-}o\text{-}tok \rightarrow kr\underline{e}v \quad tok$ bloodflow.nou blood.nou flow.nou

 $modr-o-ok\acute{y}
ightarrow modr\acute{y}
ightarrow oko$ blue-eyed.NOU blue.NOU eye.NOU

 $ps\text{-}o\text{-}v\underline{o}d o p\underline{e}s v\underline{e}st$ dog handler.NOU dog.NOU lead.VERB

 $hol\underline{e}kchtiv\hat{y} \rightarrow chtiv\hat{y}$ holka (holek = gen. pl.) wanting girls.ADJ wanting.ADJ girl.NOU

 $dvaap \'ullet \'y
ightarrow dv \check{e} \ a \ p \'u l \ l \acute{e} to$ two-and-a-half-year-old.ADJ two.NUM and.CONJ half.NUM year.AD.

```
soci-o-logie \rightarrow -soci- -log-
sociology.NOU -soci-.NEOCON -log-.NEOCON
```

 $tetrachlorethylen \rightarrow -tetra-$ chlor ethylentetrachlorethylene.NOU -tetra-.NEOCON chlorine.NOU ethylene.NOU

What makes the tasks difficult?

krv-o-tok $\rightarrow krev$ tok bloodflow.NOU blood.NOU flow.NOU $modr-o-ok\acute{y}
ightarrow modr\acute{y}
ightarrow oko$ blue-eyed.NOU blue.NOU eye.NOU

What makes the tasks difficult?

krv-o-tok $\rightarrow krev$ tok bloodflow.NOU blood.NOU flow.NOU $modr-o-ok\acute{y}
ightarrow modr\acute{y}
ightarrow oko$ blue-eved.NOU blue.NOU eye.NOU $ps ext{-}o ext{-}vod agence pes ext{ v}\acute{e}st$ dog handler.NOU dog.NOU lead.VERB

What makes the tasks difficult?

krv-o-tok $\rightarrow krev$ tok bloodflow.NOU blood.NOU flow.NOU $modr-o-ok\acute{y}
ightarrow modr\acute{y}
ightarrow oko$ blue-eyed.NOU blue.NOU eye.NOU $ps ext{-}o ext{-}vod agence pes ext{ v}\acute{e}st$ dog handler.NOU dog.NOU lead.VERB $holekchtivý \rightarrow chtivý holka$ (holek = gen. pl.) wanting girls.ADJ wanting.ADJ girl.NOU

What makes the tasks difficult?

krv-o-tok $\rightarrow krev$ tok bloodflow.NOU blood.NOU flow.NOU $modr-o-ok\acute{y}
ightarrow modr\acute{y}
ightarrow oko$ blue-eyed.NOU blue.NOU eye.NOU $ps ext{-}o ext{-}vod agence pes ext{ v}\acute{e}st$ dog handler.NOU dog.NOU lead.VERB $holekchtivý \rightarrow chtivý holka$ (holek = gen. pl.) wanting girls.ADJ wanting.ADJ girl.NOU dvaapůlletý $ightarrow dv \check{e} ~~a~~p
ull l$ léto two-and-a-half-year-old.ADJ two.NUM and.CONJ half.NUM vear.ADJ

What makes the tasks difficult?

krv-o-tok $\rightarrow krev$ tok bloodflow.NOU blood.NOU flow.NOU $modr-o-ok\acute{y}
ightarrow modr\acute{y}
ightarrow oko$ blue-eyed.NOU blue.NOU eye.NOU $ps ext{-}o ext{-}vod agence pes ext{ v}\acute{e}st$ dog handler.NOU dog.NOU lead.VERB $holekchtivý \rightarrow chtivý holka$ (holek = gen. pl.) wanting girls.ADJ wanting.ADJ girl.NOU dvaapůlletý $ightarrow dv \check{e} ~~a~~p
ull l$ léto two-and-a-half-vear-old.ADJ two.NUM and.CONJ half.NUM vear.ADJ soci-o-logie
ightarrow -soci- -logsociology.NOU -soci-.NEOCON -log-.NEOCON

5/18

What makes the tasks difficult?

krv-o-tok $\rightarrow krev$ tok bloodflow.NOU blood.NOU flow.NOU $modr-o-ok\acute{y}
ightarrow modr\acute{y}
ightarrow oko$ blue-eyed.NOU blue.NOU eye.NOU $ps ext{-}o ext{-}vod agence pes ext{ v}\acute{est}$ dog handler.NOU dog.NOU lead.VERB $holekchtivý \rightarrow chtivý holka$ (holek = gen. pl.) wanting girls.ADJ wanting.ADJ girl.NOU dvaapůlletýightarrow dv ec e ~~a~~~p ec u lléto two-and-a-half-vear-old.ADJ two.NUM and.CONJ half.NUM vear.ADJ soci-o-logie
ightarrow -soci- -logsociology.NOU -soci-.NEOCON -log-.NEOCON $tetrachlorethylen \rightarrow -tetra-$ chlor ethylentetrachlorethylene.NOU -tetra-.NEOCON chlorine.NOU ethylene.NOU

The ways (Czech) compounds have been handled so far

• Theoretical accounts:

- Bozděchová (1997) uses an onomasiological approach
- Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
- Ološtiak and Vojteková (2021) cover neoclassical compounding

• Static data resources:

- DerivAnCze: contains no compounds
- DeriNet 2.0: contained 33,923 identified and 1,252 split compounds
- Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%

The ways (Czech) compounds have been handled so far

• Theoretical accounts:

- o Bozděchová (1997) uses an onomasiological approach
- Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
- Ološtiak and Vojteková (2021) cover neoclassical compounding

• Static data resources:

- DerivAnCze: contains no compounds
- DeriNet 2.0: contained 33,923 identified and 1,252 split compounds

• Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%

The ways (Czech) compounds have been handled so far

- Theoretical accounts:
 - o Bozděchová (1997) uses an onomasiological approach
 - Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
 - Ološtiak and Vojteková (2021) cover neoclassical compounding
- Static data resources:
 - DerivAnCze: contains no compounds
 - DeriNet 2.0: contained 33,923 identified and 1,252 split compounds
- Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%

The ways (Czech) compounds have been handled so far

- Theoretical accounts:
 - o Bozděchová (1997) uses an onomasiological approach
 - Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
 - o Ološtiak and Vojteková (2021) cover neoclassical compounding
- Static data resources:
 - DerivAnCze: contains no compounds
 - DeriNet 2.0: contained 33,923 identified and 1,252 split compounds
- Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%

The ways (Czech) compounds have been handled so far

- Theoretical accounts:
 - o Bozděchová (1997) uses an onomasiological approach
 - Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
 - o Ološtiak and Vojteková (2021) cover neoclassical compounding

• Static data resources:

- DerivAnCze: contains no compounds
- DeriNet 2.0: contained 33,923 identified and 1,252 split compounds

• Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%
- Theoretical accounts:
 - o Bozděchová (1997) uses an onomasiological approach
 - Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
 - o Ološtiak and Vojteková (2021) cover neoclassical compounding
- Static data resources:
 - DerivAnCze: contains no compounds
 - DeriNet 2.0: contained 33,923 identified and 1,252 split compounds
- Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%

- Theoretical accounts:
 - o Bozděchová (1997) uses an onomasiological approach
 - Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
 - · Ološtiak and Vojteková (2021) cover neoclassical compounding
- Static data resources:
 - DerivAnCze: contains no compounds
 - DeriNet 2.0: contained 33,923 identified and 1,252 split compounds
- Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%

- Theoretical accounts:
 - o Bozděchová (1997) uses an onomasiological approach
 - Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
 - · Ološtiak and Vojteková (2021) cover neoclassical compounding
- Static data resources:
 - DerivAnCze: contains no compounds
 - $\circ~$ DeriNet 2.0: contained 33,923 identified and 1,252 split compounds
- Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%

- Theoretical accounts:
 - o Bozděchová (1997) uses an onomasiological approach
 - Štichauer (2003) discusses Bisetto and Scalise's (2005) approach wrt Czech
 - · Ološtiak and Vojteková (2021) cover neoclassical compounding
- Static data resources:
 - DerivAnCze: contains no compounds
 - DeriNet 2.0: contained 33,923 identified and 1,252 split compounds
- Procedural tools:

Authors	Language	Task	Performance
Krotova et al. (2020)	German	Splitting	Acc: 96%
Hellwig and Nehrdich (2018)	Sanskrit	Splitting	Acc: 96%
Clouet and Daille (2014)	English	Finding split-points	F1: 80%
Clouet and Daille (2014)	Russian	Finding split-points	F1: 63%

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

ryb-o-lov ightarrow ryba lov fishery.NOU fish.NOU hunt.NOU

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

ryb-o-lov $\rightarrow ryba$ lov fishery.NOU fish.NOU hunt.NO

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

ryb-o-lov $\rightarrow ryba$ lov fishery.NOU fish.NOU hunt.NO

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

ryb-o-lov
ightarrow ryba lov fishery.NOU fish.NOU hunt.NOU

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

 $ryb{-}o{-}lov
ightarrow ryba lov$ fishery.NOU fish.NOU hunt.NOU

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

ryb-o-lov
ightarrow ryba lov fishery.NOU fish.NOU hunt.NOU

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

 $ryb\text{-}o\text{-}lov \rightarrow ryba \quad lov$ fishery.NOU fish.NOU hunt.NOU

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

 $ryb\text{-}o\text{-}lov \rightarrow ryba \quad lov$ fishery.NOU fish.NOU hunt.NOU

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

 $ryb\text{-}o\text{-}lov \rightarrow ryba \quad lov$ fishery.NOU fish.NOU hunt.NOU

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

 $ryb\text{-}o\text{-}lov \rightarrow ryba \quad lov$ fishery.NOU fish.NOU hunt.NOU

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

What we had to work with

Data set: Manually annotated data

- 1,447 manually annotated compounds (by us)
- extracted from DeriNet
- 80/10/10 split (1,158/145/144; training/testing/validation)

 $ryb\text{-}o\text{-}lov \rightarrow ryba \quad lov$ fishery.NOU fish.NOU hunt.NOU

Synthetic data

- 280,000 compounds synthesized using lexemes from DeriNet
- Used only for Czech Compound Splitter training

How we evaluated performance

- Accuracy: Ratio of correct predictions vs. incorrect predictions
 - Applicable to both compound splitting and identification
 - Within splitting, a prediction is correct iff *all* parents are correct (unless stated otherwise)
- **Root accuracy**: Ratio of predicted parents within the correct derivational family
 - Only applicable to compound splitting
 - Checked using DeriNet

How we evaluated performance

- Accuracy: Ratio of correct predictions vs. incorrect predictions
 - Applicable to both compound splitting and identification
 - Within splitting, a prediction is correct iff *all* parents are correct (unless stated otherwise)
- **Root accuracy**: Ratio of predicted parents within the correct derivational family
 - Only applicable to compound splitting
 - Checked using DeriNet

How we evaluated performance

- Accuracy: Ratio of correct predictions vs. incorrect predictions
 - $\circ~$ Applicable to both compound splitting and identification
 - Within splitting, a prediction is correct iff *all* parents are correct (unless stated otherwise)
- **Root accuracy**: Ratio of predicted parents within the correct derivational family
 - Only applicable to compound splitting
 - Checked using DeriNet

How we evaluated performance

- Accuracy: Ratio of correct predictions vs. incorrect predictions
 - Applicable to both compound splitting and identification
 - Within splitting, a prediction is correct iff *all* parents are correct (unless stated otherwise)
- **Root accuracy**: Ratio of predicted parents within the correct derivational family
 - Only applicable to compound splitting
 - Checked using DeriNet

How we evaluated performance

- Accuracy: Ratio of correct predictions vs. incorrect predictions
 - Applicable to both compound splitting and identification
 - Within splitting, a prediction is correct iff *all* parents are correct (unless stated otherwise)
- **Root accuracy**: Ratio of predicted parents within the correct derivational family
 - Only applicable to compound splitting
 - Checked using DeriNet

How we evaluated performance

- Accuracy: Ratio of correct predictions vs. incorrect predictions
 - Applicable to both compound splitting and identification
 - Within splitting, a prediction is correct iff *all* parents are correct (unless stated otherwise)
- **Root accuracy**: Ratio of predicted parents within the correct derivational family
 - Only applicable to compound splitting
 - Checked using DeriNet

How we evaluated performance

- Accuracy: Ratio of correct predictions vs. incorrect predictions
 - Applicable to both compound splitting and identification
 - Within splitting, a prediction is correct iff *all* parents are correct (unless stated otherwise)
- **Root accuracy**: Ratio of predicted parents within the correct derivational family
 - Only applicable to compound splitting
 - Checked using DeriNet

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

Three approaches to the tasks at hand

• Baseline solution

- Naive explicit solution contextualizing the other two
- Limited to compound splitting
- Returns two parents
- Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - $\circ~$ Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - $\circ~$ Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - $\circ~$ Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - $\circ~$ Deep learning solution
 - $\circ~$ Performs both compound identification and splitting
 - $\circ~$ Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon
Our solutions

Three approaches to the tasks at hand

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

Our solutions

Three approaches to the tasks at hand

- Baseline solution
 - Naive explicit solution contextualizing the other two
 - Limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- IML()-based heuristic algorithm
 - Explicit solution based on a phonological similarity function (IML())
 - Also limited to compound splitting
 - Returns two parents
 - Requires a lexicon
- Czech Compound Splitter
 - Deep learning solution
 - Performs both compound identification and splitting
 - Returns any number of parents
 - Does not require a lexicon

The *IML*() (*Interlexical Matrices of Likeness*) function returns the phonological similarity of two Czech graphical words using a manually defined *correspondence matrix*.

- $IML(word_1, word_2) = 0$ iff $word_1 = word_2$
- *IML*(*word*₁, *word*₂) need not equal *IML*(*word*₂, *word*₁)

The IML()-based heuristic algorithm tests all word-pairs ((*lexeme*₁, *lexeme*₂)) from the lexicon that satisfy a heuristic filter that discards all obviously wrong word-pairs.

• otherwise, $800,000^2 = 6.4 \times 10^{11}$ operations

The *IML*() (*Interlexical Matrices of Likeness*) function returns the phonological similarity of two Czech graphical words using a manually defined *correspondence matrix*.

- $IML(word_1, word_2) = 0$ iff $word_1 = word_2$
- *IML*(*word*₁, *word*₂) need not equal *IML*(*word*₂, *word*₁)

The IML()-based heuristic algorithm tests all word-pairs ((*lexeme*₁, *lexeme*₂)) from the lexicon that satisfy a heuristic filter that discards all obviously wrong word-pairs.

- otherwise, $800,000^2 = 6.4 \times 10^{11}$ operations

The *IML*() (*Interlexical Matrices of Likeness*) function returns the phonological similarity of two Czech graphical words using a manually defined *correspondence matrix*.

- $IML(word_1, word_2) = 0$ iff $word_1 = word_2$
- *IML*(*word*₁, *word*₂) need not equal *IML*(*word*₂, *word*₁)

The IML()-based heuristic algorithm tests all word-pairs ((*lexeme*1, *lexeme*2)) from the lexicon that satisfy a heuristic filter that discards all obviously wrong word-pairs.

• otherwise, $800,000^2 = 6.4 \times 10^{11}$ operations

The *IML*() (*Interlexical Matrices of Likeness*) function returns the phonological similarity of two Czech graphical words using a manually defined *correspondence matrix*.

- $IML(word_1, word_2) = 0$ iff $word_1 = word_2$
- *IML*(*word*₁, *word*₂) need not equal *IML*(*word*₂, *word*₁)

The IML()-based heuristic algorithm tests all word-pairs ((*lexeme*1, *lexeme*2)) from the lexicon that satisfy a heuristic filter that discards all obviously wrong word-pairs.

• otherwise, $800,000^2 = 6.4 \times 10^{11}$ operations

The *IML*() (*Interlexical Matrices of Likeness*) function returns the phonological similarity of two Czech graphical words using a manually defined *correspondence matrix*.

- $IML(word_1, word_2) = 0$ iff $word_1 = word_2$
- *IML*(*word*₁, *word*₂) need not equal *IML*(*word*₂, *word*₁)

The IML()-based heuristic algorithm tests all word-pairs $((lexeme_1, lexeme_2))$ from the lexicon that satisfy a heuristic filter that discards all obviously wrong word-pairs.

- otherwise, $800,000^2 = 6.4 \times 10^{11}$ operations

The *IML*() (*Interlexical Matrices of Likeness*) function returns the phonological similarity of two Czech graphical words using a manually defined *correspondence matrix*.

- $IML(word_1, word_2) = 0$ iff $word_1 = word_2$
- *IML*(*word*₁, *word*₂) need not equal *IML*(*word*₂, *word*₁)

The IML()-based heuristic algorithm tests all word-pairs ((*lexeme*₁, *lexeme*₂)) from the lexicon that satisfy a heuristic filter that discards all obviously wrong word-pairs.

- otherwise, $800,000^2 = 6.4 \times 10^{11}$ operations

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

It was trained on:

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

 $egin{array}{ccc} egin{array}{ccc} egin{array}{cccc} egin{array}{ccc} egin{array}{ccc} egin{arr$

Czech Compounds Splitter (*CCS*) performs compound splitting and compound identification in one step.

It is an implementation of the Marian translator framework by Microsoft.

It was trained on:

- 1,164 genuine compounds, with their splittings
- 280,000 synthetic compounds, with their splittings
- the near entirety of DeriNet's non-compounds, with their derivational parents

 $egin{array}{ccc} egin{array}{ccc} egin{array}{cccc} egin{array}{ccc} egin{array}{ccc} egin{arr$

Czech Compounds Splitter's typical errors in compound splitting:

Compound	CCS splitting	Correct splitting		
dlouhohořící	dlouhohořící	dlouho hořící		
CCS returns t	he original string	, performing no splitting.		
osmiramenný	osm ramenný	osm rameno		
CCS returns a non-existing derivative of an existing word.				
bíločerný	bílo černý	bílý černý		
CCS includes	the interfix in on	e of the parents.		

Czech Compounds Splitter's typical errors in compound splitting:

Compound	CCS splitting	Correct splitting		
dlouhohořící	dlouhohořící	dlouho hořící		
CCS returns t	he original string	, performing no splitting.		
osmiramenný	osm ramenný	osm rameno		
CCS returns a non-existing derivative of an existing word.				
bíločerný	bílo černý	bílý černý		
CCS includes	the interfix in on	e of the parents.		

Czech Compounds Splitter's typical errors in compound splitting:

Compound	CCS splitting	Correct splitting		
dlouhohořící	dlouhohořící	dlouho hořící		
CCS returns the original string, performing no splitting.				
osmiramenný osm ramenný osm rameno				
CCS returns a non-existing derivative of an existing word.				
bíločerný	bílo černý	bílý černý		
CCS includes the interfix in one of the parents.				

Results

How the approaches ended up performing

Compound identification:

- Performed only by Czech Compound Splitter
- Evaluated on 144 compounds and 141 non-compounds (no information rate: 50.5%)
- Accuracy: 92%

Compound splitting:

	1st parent	2nd parent	Overall	Overall root
Method	accuracy	accuracy	accuracy	accuracy
Baseline	22%	42%	11%	16%
IML()	42%	66%	24%	39%
CCS	61%	66%	54%	61%

- We have made the first attempt into the automatic splitting and identification of Czech compound words.
- We have developed three tools:
 - a naive baseline solution,
 - a heuristic algorithm based on IML(), a custom phonological similarity function
 - Czech Compound Splitter, a deep learning solution
- *Czech Compound Splitter* turns out to be the best solution, in terms of both performance and practicality.
- It has achieved an accuracy of 54% in compound splitting and an accuracy of 92% in compound splitting.
- Furthermore, it can also find the single *derivational* parent of a given Czech word, if it's not a compound.

- We have made the first attempt into the automatic splitting and identification of Czech compound words.
- We have developed three tools:
 - $\circ\,$ a naive baseline solution,
 - a heuristic algorithm based on IML(), a custom phonological similarity function
 - Czech Compound Splitter, a deep learning solution
- *Czech Compound Splitter* turns out to be the best solution, in terms of both performance and practicality.
- It has achieved an accuracy of 54% in compound splitting and an accuracy of 92% in compound splitting.
- Furthermore, it can also find the single *derivational* parent of a given Czech word, if it's not a compound.

- We have made the first attempt into the automatic splitting and identification of Czech compound words.
- We have developed three tools:
 - a naive baseline solution,
 - a heuristic algorithm based on IML(), a custom phonological similarity function
 - Czech Compound Splitter, a deep learning solution
- *Czech Compound Splitter* turns out to be the best solution, in terms of both performance and practicality.
- It has achieved an accuracy of 54% in compound splitting and an accuracy of 92% in compound splitting.
- Furthermore, it can also find the single *derivational* parent of a given Czech word, if it's not a compound.

- We have made the first attempt into the automatic splitting and identification of Czech compound words.
- We have developed three tools:
 - a naive baseline solution,
 - a heuristic algorithm based on IML(), a custom phonological similarity function
 - Czech Compound Splitter, a deep learning solution
- *Czech Compound Splitter* turns out to be the best solution, in terms of both performance and practicality.
- It has achieved an accuracy of 54% in compound splitting and an accuracy of 92% in compound splitting.
- Furthermore, it can also find the single *derivational* parent of a given Czech word, if it's not a compound.

- We have made the first attempt into the automatic splitting and identification of Czech compound words.
- We have developed three tools:
 - a naive baseline solution,
 - a heuristic algorithm based on IML(), a custom phonological similarity function
 - Czech Compound Splitter, a deep learning solution
- *Czech Compound Splitter* turns out to be the best solution, in terms of both performance and practicality.
- It has achieved an accuracy of 54% in compound splitting and an accuracy of 92% in compound splitting.
- Furthermore, it can also find the single *derivational* parent of a given Czech word, if it's not a compound.

- We have made the first attempt into the automatic splitting and identification of Czech compound words.
- We have developed three tools:
 - a naive baseline solution,
 - a heuristic algorithm based on IML(), a custom phonological similarity function
 - Czech Compound Splitter, a deep learning solution
- *Czech Compound Splitter* turns out to be the best solution, in terms of both performance and practicality.
- It has achieved an accuracy of 54% in compound splitting and an accuracy of 92% in compound splitting.
- Furthermore, it can also find the single *derivational* parent of a given Czech word, if it's not a compound.

- We have made the first attempt into the automatic splitting and identification of Czech compound words.
- We have developed three tools:
 - a naive baseline solution,
 - a heuristic algorithm based on IML(), a custom phonological similarity function
 - Czech Compound Splitter, a deep learning solution
- *Czech Compound Splitter* turns out to be the best solution, in terms of both performance and practicality.
- It has achieved an accuracy of 54% in compound splitting and an accuracy of 92% in compound splitting.
- Furthermore, it can also find the single *derivational* parent of a given Czech word, if it's not a compound.

- We have made the first attempt into the automatic splitting and identification of Czech compound words.
- We have developed three tools:
 - a naive baseline solution,
 - a heuristic algorithm based on IML(), a custom phonological similarity function
 - Czech Compound Splitter, a deep learning solution
- *Czech Compound Splitter* turns out to be the best solution, in terms of both performance and practicality.
- It has achieved an accuracy of 54% in compound splitting and an accuracy of 92% in compound splitting.
- Furthermore, it can also find the single *derivational* parent of a given Czech word, if it's not a compound.

This work was supported by the Grant No. GA19-14534S of the Czech Science Foundation, the LINDAT/CLARIAH-CZ project of the Ministry of Education, Youth and Sports of the Czech Republic (project LM2018101), and by the Grant No. START/HUM/010 of Grant schemes at Charles University (Reg. No. CZ.02.2.69/0.0/0.0/19_073/0016935).

References



http://hdl.handle.net/11234/1-3765.

References 2

GermaNet - a lexical-semantic net for German.

Hamp, B and Feldweg, H. In: Automatic information extraction and building of lexical semantic resources for NLP applications, 1997.



The CELEX lexical database.

Baayen, H, Piepenbrock, R, and Gulikers, L. Linguistic Data Consortium, 1996.



A Joint Approach to Compound Splitting and Idiomatic Compound Detection.

Krotova, I, Aksenov, S, Artemova E. In: *Proceedings of the 12th Language Resources and Evaluation Conference*, 2020.

Sanskrit word segmentation using character-level recurrent and convolutional neural networks. Hellwig, O, Nehrdich, S. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018.

Splitting of compound terms in non-prototypical compounding languages.

Clouet, E. and Daille, B. In: Workshop on Computational Approaches to Compound Analysis, 2014

Thank you for your attention!

Emil Svoboda

svoboda@ufal.mff.cuni.cz