


HiČKoK: History of Czech in Corpus Continuum

Daniel Zeman

 December 9, 2024



Three-year project (September 2023 – November 2026)



Funded by TAČR

Carried out by



Národní knihovna
České republiky
National Library
of the Czech Republic



Three-year project (September 2023 – November 2026)



Funded by TAČR

Carried out by



- 1 Goals
- 2 Old Czech: Workflow
- 3 Old Czech: Interesting Issues

Goals

Výběr řádků: základní

1 / 1 347

zalokovaná, za což jsi velice vděčný. </p><p> A hopla jednoho jsou tak akorát pro hlodavce, na chodbě se to máš nechat plavat. Zatímco jsi byl pryč, měla barvy pšenice, aspoň tak sis to představoval; , že by sis ji mohl s Dolly spojit. </p><p> je svým lehkomyšlným způsobem velkorysý. " Máte k sobě tam našla chyby? " </p><p> Má seznam pěkně po ruce Alex Hardy. Vážně pokývne hlavou, jako byste právě peníze z pojistky - jistě je jen, že zmizely stříbrnými vlasy pečlivě zkoumá tvou pozvánku. Po stranách dveří paměti a pátráš po nějakých stopách. </p><p> Ve vrchním šuplíku . Kreditní kartou ho vyškrábneš na stůl a kartou také , že má rozkošný zadeček. Dělal jsi s ní je velká a lesklá a je na ní mladší Megan ten česnek, " zeptá se. Podářilo se ti radši natáhl. Zabořit se hlavou Megan do klína a že věkem bude ještě pod zákonem. Pod očima si namalovala/namalovat/VpFS---R-AA---P dva/dva/CIIP4----- namalovala/namalovat/VpFS---R-AA---P dva/dva/CIIP4-----

,/;/Z:----- dva/dva/CIIP1----- ,/;/Z:-----
 vyhnou/vyhnut/VB-P---3P-AA---P dva/dva/CIIP1----- chodci/chodec/NNMP1-----A-----
 brala/brát/VpFS---R-AA---I dva/dva/CIIP4----- vzkazy/vzkaz/NNIP4-----A-----
 za/za/RR--4----- dva/dva/CIIP4----- měsíce/měsíc/NNIP4-----A-----
 O/o/RR--4----- dva/dva/CIIP4----- roky/rok/NNIP4-----A-----
 vy/vy/PP-P1--2----- dva/dva/CIIP1----- hodně/hodně/Dg-----1A-----
 :;/Z:----- dva/dva/CIIP4----- opačné/opačný/AAIP4---1A-----
 vy/vy/PP-P1--2----- dva/dva/CIIP4----- byli/být/VpMP----R-AA---I
 za/za/RR--4----- dva/dva/CIIP4----- týdny/týden/NNIP4-----A-----
 stojí/stát/VB-P---3P-AA---B dva/dva/CIIP1----- ohromní/ohromný/AAMP1---1A-----
 objevíš/objevit/VB-S---2P-AA---P dva/dva/CIIP4----- prázdné/prázdný/AAIP4---1A-----
 uděláš/udělat/VB-S---2P-AA---P dva/dva/CIIP4----- rovné/rovný/AAFP4---1A-----
 skoro/skoro/Db----- dva/dva/CIIP4----- roky/rok/NNIP4-----A-----
 a/a/J^----- dva/dva/CIIP1----- muži/muž/NNMP1-----A-----
 oloupat/oloupat/Vf-----A---P dva/dva/CIIP4----- stroužky/stroužek/NNIP4-----A-----
 týden/týden/NNIS4---A--- dva/dva/CIIP1----- tak/tak/Db-----
 namalovala/namalovat/VpFS---R-AA---P dva/dva/CIIP4----- purpurové/purpurový/AAIP4---1A-----
 Uděláš/udělat/VB-S---2P-AA---P dva/dva/CIIP4----- řádky/řádek/NNIP4-----A-----
 bejt/být/Vf-----A---6I dva/dva/CIIP1----- večery/večeř/NNIP1-----A-----
 První/první/CrMP1----- dva/dva/CIIP1----- jsou/být/VB-P---3P-AA---I
 ti/ty/PH-S3--2----- dva/dva/CIIP4----- tlusté/tlustý/AAIP4---1A-----

tři, čtyři . Vojáci už se postavili na nohy . </p><p> Pluješ po linoleu do Oddělení verifikace , jeden od jakéhosi monsieur a z oddělení čehosi a v Kansasu jsi ještě žádnou pšenici neviděl . Většinu později byla Dolly pozvaná na svatbu na východě . blízko, " chceš vědět . </p><p> " Podle mě je přízvuky, volební obvod ve střední Francii nesprávě osobami, které hodlá najít o půlnoci v Kláfině kanceláři . </p><p> U vchodu vysoká žena se stříbrnými vlasy černošlákem a paže mají zkřížené na prsou obdélíkové balíčky . Jeden z nich vlastně úplně pro řádky . Podíváš se na Megan . Čte . Klidně a vůbec sis té nádhery nevšímá . Kolik jí vlastně na jevišti . </p><p> " To byla moje poslední role . Jsou pěkně nahatě . </p><p> " Moc výkonní zrovna zůstat . Postel je jen pár kroků odsud . Nakloníš hlavu a připomínat lícní kosti . Už víš a sedneš si do křesla . Dnes je to přesně za sebou a samozřejmě, že expres Allagash vykoleno plně kokainových fandů a seriózních žvanilů . Třetí krajice . Kůrku měly spálenou, ale uvnitř byly navlhčené, které mám obzvlášť rád . " Velikou, tajemnou : pošlapání Spravedlnosti a vláda Tyrannie . Ty ostatní dojem: hlučný sobotní trh, rozpínající se po městě , filonjouovský, který navrhoval francouzského kavaléra a stále ještě nebylo vidět světlo habemus nad námi

▲ Knoflík na dveřích se otáčí sem a tam . Allagash tě šfouchá do ruky a rty se horečně vpytává . Zámek cvakne a dveře se rozletí . Stojí v nich Alex Hardy . Vážně pokývne hlavou , jako byste právě vy dva byli osobami , které hodlá najít o půlnoci v Kláfině kanceláři . Snažíš se rychle vymyslet povídačku , kterou by spolkl . Tad třímá dřevěný metr , který našel za dveřmi . </p><p> " Ty jsi nás , Alexi , vylekal . ▼

Výběr řádků: základní

1 / 12 470

Inga Ābele Ostřice Hra o zmlkne .) A tak dál a tak dál .	dvou /NUM/Case=Loc Number=Plur NumForm=Word NumType=Card	dějstvích Osoby Vilhelms Putns Lija , jeho že
vináf ke mně přišel s lahví koňaku , že jsem . Lija . Aleksisi , přečti ještě něco . Ti	Dvě /NUM/Case=Nom Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	stránky a fotografie . Všichni to čtou , celé m
Vilhelms . Teď je to snad v pořádku . Už děje jenom v naší čtvrti . Lija . Měla jsem	dva /NUM/Animacy=Inan Case=Acc Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	dny nejedla ... , hodně toho tam není napsán
Aleksis . To je teda vážně ostuda . A já knem do tmy . Hledá cigarety , ale najde jen	dva /NUM/Animacy=Anim Case=Nom Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	už mě omrzeli . leva přistoupí ke stolu a sebi
... Ale neřekla jsem ti , že ve skutečnosti už si vyšfourával zbytky sena . Po jednom či po	dva /NUM/Animacy=Inan Case=Acc Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	roky je pryč . Teď si ho sem přivolala mamín
očí stisk čelistí . Na břehu ležely na dekách dcera a syn , žena , muž a sestra ,	dva /NUM/Animacy=Anim Case=Acc Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	syny – vlastního Aleksise a ještě jednoho . V
a sestra , dvě sestry s bratrem nebo otec s lo malého města , vrásčitá , unavená žena s	dvě /NUM/Case=Acc Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	roky žil jak anděl . Santa (Usměje se ,
mĺčky . Bertolds udělal jeden krok , chlapec ejme , že jsem ještě nic nepropásl . Koupím	dvou /NUM/Case=Loc Number=Plur NumForm=Word NumType=Card	prázdné krabičky . Zapíná rádio . Přeladuje s
. Spěchaly dál přes hřbitov . U vchodu , kde a sestra , dvě sestry s bratrem nebo otec s	dvě /NUM/Case=Nom Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	měsíce žiju s jiným týpkem , jo , no ,
Prohodila to jen tak , mezi řečí , možná mezi jet do Buenos Aires . , K tomu musej být	dvě /NUM/Case=Nom Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	se děti vracely zpět k autu a sypaly otcí do
ýho všimnul a ocenil ho , k tomu musej být že pozoruje každý její pohyb . Bílá skvrna se	dvě /NUM/Case=Nom Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	tělnaté ženy a vítaly oba příchozí . U nikoho :
Nechce se jim bez ní odejít . Pokládají ji na řed očima . Nazelenalá kost orámuje krajinu	dvěma /NUM/Animacy=Inan Case=Ins Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	sestry s bratrem nebo otec s dvěma dcerami
sklem auta , tento spánek vyvolal vizi o nich chalupu . Potraviny spravedlivě rozdělili do	dvěma /NUM/Animacy=Inan Case=Ins Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	dcerami . Jejich radost a výkřiky byly urputn
ela si s obavami leva a neslyšela žádné raz chaluhy .	dva /NUM/Animacy=Anim Case=Nom Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	nevychovanými syny . Pracovala jako poštač
	dva /NUM/Animacy=Anim Case=Nom Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	. V každé ruce se mu pohupoval sandál . Za
	dvě /NUM/Case=Acc Gender=Neut Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	pořádná kola , stan a vezmu si dovolenou . F
	dvě /NUM/Case=Nom Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	borovice připomínaly tajemný vstup na one
	dvěma /NUM/Animacy=Inan Case=Ins Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	sklenkami vína . Ale moje fantazie už začala
	dva /NUM/Animacy=Anim Case=Nom Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	,' přerušil ji vždycky muž , který tu teď
	dva /NUM/Animacy=Anim Case=Nom Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	. ~ # ~ ~ Regina se odmlčí , nanovo se
	dvě /NUM/Case=Acc Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	žhnoucíma , dehtově černýma očima na vzor
	dvě /NUM/Case=Acc Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	suché smrkové větve tak , aby se očními otv
	dvě /NUM/Case=Acc Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	okny . Jsou tak široká , že se musí dívat
	dvě /NUM/Case=Acc Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	, o ní a Pávilsovi v prvotním lese , ukázal
	dvě /NUM/Case=Acc Gender=Fem Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	baťohů , jen dočerna vyuzený kančí hřbet vze
	dva /NUM/Animacy=Anim Case=Nom Gender=Masc Number=Plur NumForm=Word NumType=Card NumValue=1 NumValue=2 NumValue=3	tři . Ve světle reflektorů viděla Pávilsu pro

Výchozí zobrazení | Formátovaný text

▲ dobře možné , že se ze slona narodil První Sníh ! Dino ! Vlk vyjde ven . Dvěje jsou otevřené . Silně prší . Po schodech schází Aleksis a vyhlíží oknem do tmy . Hledá cigarety , ale najde jen **dvě** prázdné krabičky . Zapíná rádio . Přeladuje stanice . Vypíná rádio . Opravuje kresbu . Sedne si ke stolu . Lehne si na pohovku . Přistoupí k oknu , protože přijíždí auto . Reflektory ho na moment osvětlí a auto ▼



Výběr řádků: základní ▾

1 / 33 ▶▶▶

<input type="checkbox"/>	1492	mi se , že jest . V tom městě byli	sme dva dní	, chleba vody i rozličného ovoce bylo co dosti kúpiti
<input type="checkbox"/>	1492	i na vodu , neb do něho teče sedmdesát potokův	a dva .	Pak na kupecké věci jest převelmi bohaté , ješto na
<input type="checkbox"/>	1492	. Pak opět dne jiného ráno vstavše , šli se	mnú dva bosáci	do Jeruzaléma , neb hora Sion jest od Jeruzaléma jako
<input type="checkbox"/>	1492	kus zdi kostelnie , a tak mi pravili , že	zabilo dva černé	křesťany , a nedadie jim toho zase opravití . V
<input type="checkbox"/>	1492	ješto držie křest , svétie toliko neděli , ti mají	jedno dva štrychy	nebo znamenie na tvári . Třetí toliko držie křest ,
<input type="checkbox"/>	1492	a některé muož tluščie býti . Vinic nekopají , ale	zapřaha dva buvoly	do pluhu i obvorávají je na vsecky strany , totižto
<input type="checkbox"/>	1492	vašeho . V Gaze byli sme týden . Odtud sme	jeli dva dní	i přijeli sme k jednomu městečku a neviděli sme ho
<input type="checkbox"/>	1492	daktylův jako chmele , že jich z toho kolečka natrhá	se dva věrtele	a z druhého viec . A na jednom dřevě bude
<input type="checkbox"/>	1492	toho voda z trub teče . A u toho koryta	sedíta dva starci	a ktož koli přijde k té mřeži a chce píti
<input type="checkbox"/>	1492	, a když sem toho pilen byl , potom mi	poručil dva a	tak dále , až mi i jiné věci poručel .
<input type="checkbox"/>	1492	byla , tehdy na osle jede a okolo nie jdú	služebníci dva nebo	třie a tvář jejé jest vsecka zakryta čistým tafartem černým
<input type="checkbox"/>	1492	a netoliko já , ale všichni kupcí . Byla se	nám dva pacholky	roznemohla , kteréžto tajně mezi seub chovali sme a nemocné
<input type="checkbox"/>	1835	ještě vypravoval , kdyby ho přichozí nebyli vytrhli . Byli	to dva cikáni	. Napřed šel malý , po maďarsku vystrojený mladík ,
<input type="checkbox"/>	1835	Mladý cikán postoupil , spatřiv vysloužence a ostatní , zaražený	as dva kroky	zpět , jako by byl někoho jiného očekával ; opět
<input type="checkbox"/>	1835	mezi křovím nějakou rozsedinou dolů . Starší lezl za ním	asi dva sáhy	hluboko , a oba stanuli v skalnaté kotlině asi dva
<input type="checkbox"/>	1835	dva sáhy hluboko , a oba stanuli v skalnaté kotlině	asi dva sáhy	dlouhé a sáh široké . Kolem nich strměla písečná skála
<input type="checkbox"/>	1835	se vám vracím jednou domů ; jest to ode dneška	za dva měsíce	osm let , já jsem si to zapsal . -
<input type="checkbox"/>	1835	a ještě jeden šli před nimi s rozžatými pochodněmi .	Ostatní dva ,	pomohše jim mrtvol přendati přes kostničí kámen , šli ozbrojeni
<input type="checkbox"/>	1835	prsa a začlánělo hlubokou ránu jejich . Zachvěl se cikán	; dva rychlé	kroky přistoupil bliže a po jeho nad hlavu vyzvižené pravici
<input type="checkbox"/>	1835	tlačilo nesčíslné množství vůkolního lidu . Na nejvyšším jeho vrcholku	stály dva mladé	stromky , uprostřed nich strměl čekán a kolem toho se
<input type="checkbox"/>	1835	straně kněz , no druhá mistr napraví . Za ním	vedli dva pacholky	bláznivou Angelínu uvázanou na dlouhém provaze , která do výšky
<input type="checkbox"/>	1792			přišli , kterouž ona podlé
<input type="checkbox"/>	1792			zlatý Radoj , jak
<input type="checkbox"/>	1869			výtiscích ještě na skladě . Belletristický
<input type="checkbox"/>	1869			
<input type="checkbox"/>	1869			

▲ na palmě , neb když z daleka na ně hledí , zdá se jako vohánka , pak v tom větvěví sú kola vyrostlá jako melé a na tom melí bude daktylův jako chmele , že jich z toho kolečka natrhá **se dva věrtele** a z druhého viec . A na jednom dřevě bude takových kol čtyři nebo pět . Také v tom městečku zdá mi se , že největšie živnost jich jest daktyle , neb odtud dodávají daktylův do všech zemí okolních kupcí ▼





Vokabulář webový

Webové hnízdo pramenů k poznání historické češtiny



Co je VW Aktuality Slovníky staré češtiny Korpusy Edice Mluvnice Digitalizované slovníky Odborná literatura Audioknihy Zdroje Nástroje Kontakty a odkazy Připomínky

dva

Hledej

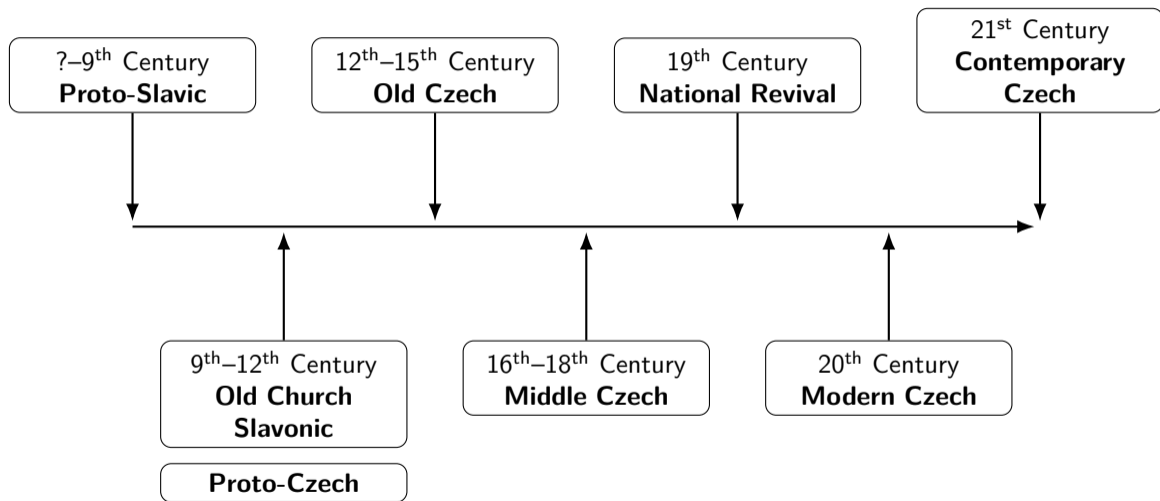
ptáci k lepu ; sto jich pobitých odvekú a cis'uc jich pod zed tekú . Takž ten šturm stá ješče více nežli pilná	dva	mésieec ; dnm i nocú bez přestánie nódpočinu šturmovánie ; jakž obirtli tako praký , nódpočinúce čas vs
řech : „ Juž ... pamět jich ... učinil toto ... “ ... múdrosti ... skonánie prozřeli ... honieše jeden tisíc a	dva	... desět tisícóc . Však proto , že böh ... prodal jě a hospodin zaklopil jě . Nenie vedě böh náš jako bohov
ověčieho těla jsa . Rovný otcu podlé božstvie , menši otcě podlé človečstvie . Jenže ač böh jest a člověk , ne	dva	však , ale jeden jest Kristus . Jeden zajisté ne obráčením božstvie v tělo , ale přijetím človečstvie u božě .
řiven , buď s tobú vet ! “ Pakliť sě lépe rozmysli , ješčeť lepši čin vymysli , jímžtoť chlapa nechudi , avšak	dva	voly vylúdi . Jako byl jeden učinil vládyka , že byl vylúdil dva voly , s hnědým plavého , u svého kmetě d
čin vymysli , jímžtoť chlapa nechudí , avšak dva voly vylúdi . Jako byl jeden učinil vládyka , že byl vylúdil	dva	voly , s hnědým plavého , u svého kmetě dobrého , jehož obviniti chtieše , ale sličně nemožieše . Pozvav j
ominaj , také mně něco z diela daj ! “ „ Nejmám , pane , nice toho , rozdal sem málo i mnoho ; kromě ješče	dva	voly jmám , oba na ofěřě prodám , na vosku a na přikrově , v němžto mě pohřebú v rově . “ „ Ješče máš d
o vás jiní . A když jima to vše povědě , již na smrtnej posteli sedě , kakž kolivěk mdle dycháše , proti diáblu	dva	dni bojováše , prosě za vše své milé syny , aby Böh spustil jich viny , i za jiné , ktož budú jeho slúhy , aby
... .. přiblíži ... ; ... jeho ne... ; ... k němu kráti ... , ... jeho všichni A když sě jako ze sna probudi ,	dva	a sedmdesát sěoci vzbudi ; smrt vsěd... na jeho hrudi , hned z něho duši vypudi . Tak sě on svých hi
edno pláčě jej zbýváše , protož ústavně plakáše ; a když sě plačie nakloni a prozřevši v hrob sě vkloni , uzře	dva	anděly z nebe u bielém rúšě blíz sebe , vniž k hlavě a k nohám více , nevzdál ot sebe sediece . Otázasta ji
ám , k dolóm i k skalám , prosiece jich , aby je skryly , na ně padnúce , přěd mě jíti nesmějice . A v tu dobu	dva	zloděje s ním vediechu , nad nimažto popravití chtiechu . Když Ježúšě vyvedú z města na popravu , tehdy
Ježiš , vecě k němu : „ Šimone , jmám tobě cos pověděti . “ A on vecě : „ Místře , pověz ! “ I vecě Ježiš : „	Dva	biešta dlužna jednomu lichevníku , jeden bieše dlužen patset peněz a druhý patdesát . A když nejmějista c
modliti . I stala sě jest krása jeho obličjě jináká , když sě modlěše , a jeho rúcho bielě stkvieše sě . A tehdy	dva	mužě mluviešta s ním . A to biešta Mojžieš a Heliáš , viděna u veleslavenství , a mluviešta o jeho přihodě
lémě . Ale Petr a ta , ješto s ním jdiešta , otrápeni biechu snem . A procitvše uzřechu jeho veleslavenství a	dva	mužě , ješto s ním stojiešta . I stalo sě jest , když otjidešta ot něho ta dva mužě , vecě Petr k Ježišěvi : „ Př
vše uzřechu jeho veleslavenství a dva mužě , ješto s ním stojiešta . I stalo sě jest , když otjidešta ot něho ta	dva	mužě , vecě Petr k Ježišěvi : „ Přikazateli , dobré jest nám tuto býti . A učiněme tři stany : jeden tobě a jede
ú ruku ku pluhu a ozřě sě za sě , nebude klčen k království božiemu . “ Pak potom znamena hospodin jiných	dva	a sedmdesát i rozesla jě po dvú do každého města vsudy , kamž jmějieše sám přijíti . I mluvieše jim : „ Jis
řší , mě slyši , a kto vámi hrdá , mnú hrdá , a kto mnú hrdá , hrdá tiem , kto mě jest poslal . “ I vrátilo sě jest	dva	a sedmdesát učedníkuov s radosti řkúce : „ Hospodine , také i běsové jsú nám poddáni ve jmě tvě . “ I vec
řší , mě slyši , a kto vámi hrdá , mnú hrdá , a kto mnú hrdá , hrdá tiem , kto mě jest poslal . “ I vrátilo sě jest	dva	paniazě , i da otrolu žla : Imái jim názi , a což kali naž položě , iáz kděš sě vrátim , vráci tobě . Kto z těj

doc Bible drážďanská, Lukášovo evangelium; 2. polovina 60. let 14. století; rukopis; Německo; Drážďany; Saská státní knihovna – Státní a univerzitní knihovna v Drážďanech (Die Sächsische Landesbibliothek – Staats- und Universitätsbibliothek Dresden); Mscr.Dresd.Oe.85; 14; Bibl; BiblDrážď; próza; bibliický text; [ediční poznámka](#); Hlaváčová Svobodová, Andrea (Ústav pro jazyk český AV ČR, v. v. i.) – Voleková, Kateřina (Ústav pro jazyk český AV ČR, v. v. i.)

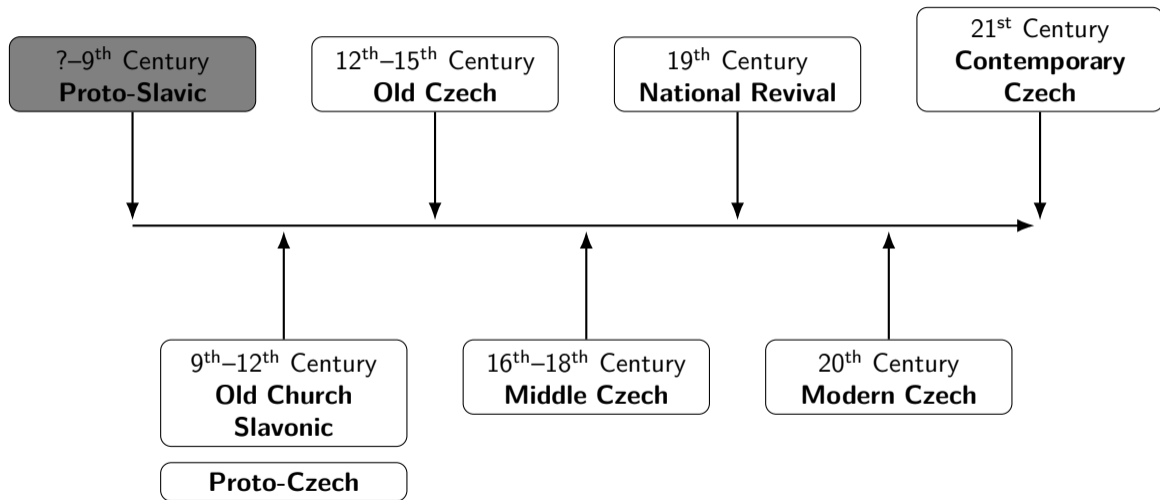
Goal

- Manual annotation of texts from all stages of Czech
- **Universal Dependencies** annotation style
- **Only morphology** (lemma, UPOS, features); no syntax
- Large enough to train a model ('etalon', ca. 100K words per period)
- UDPipe models will be freely available
 - Unfortunately, we are not allowed to make the data available for download
- Sources of data: ÚJČ, ČNK, NK
 - For some stages, no large corpus was available so far
 - Not even unannotated

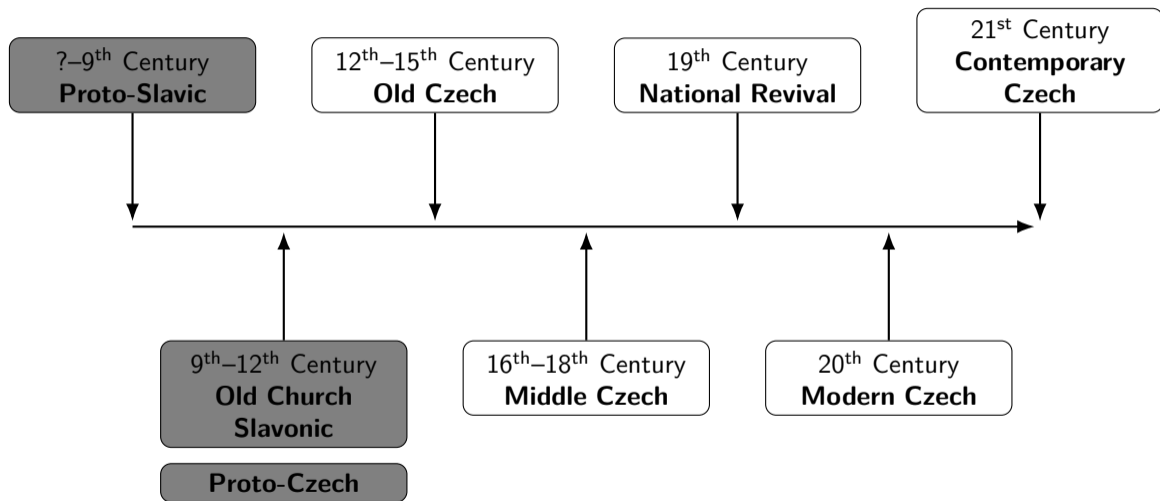
History of the Czech Language



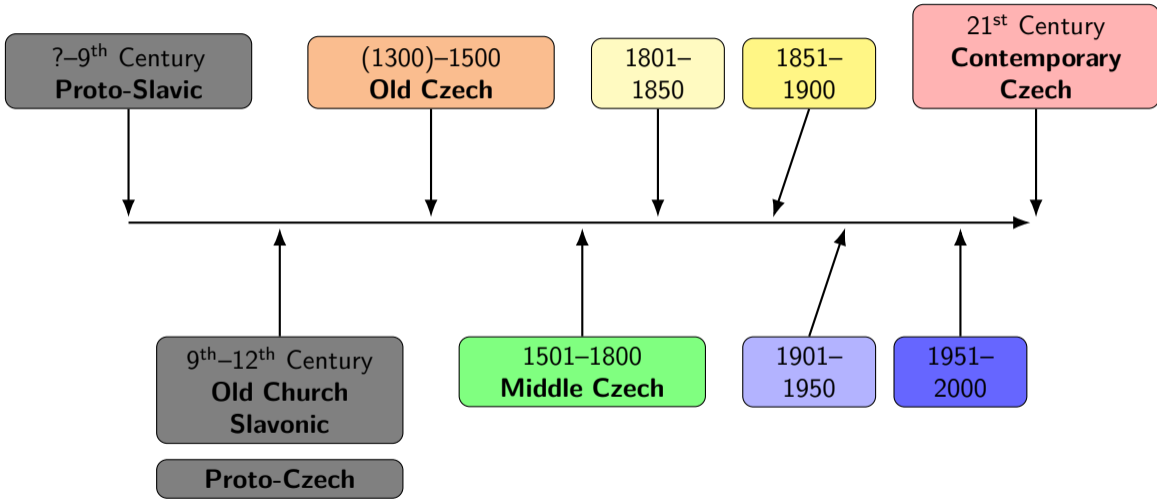
History of the Czech Language



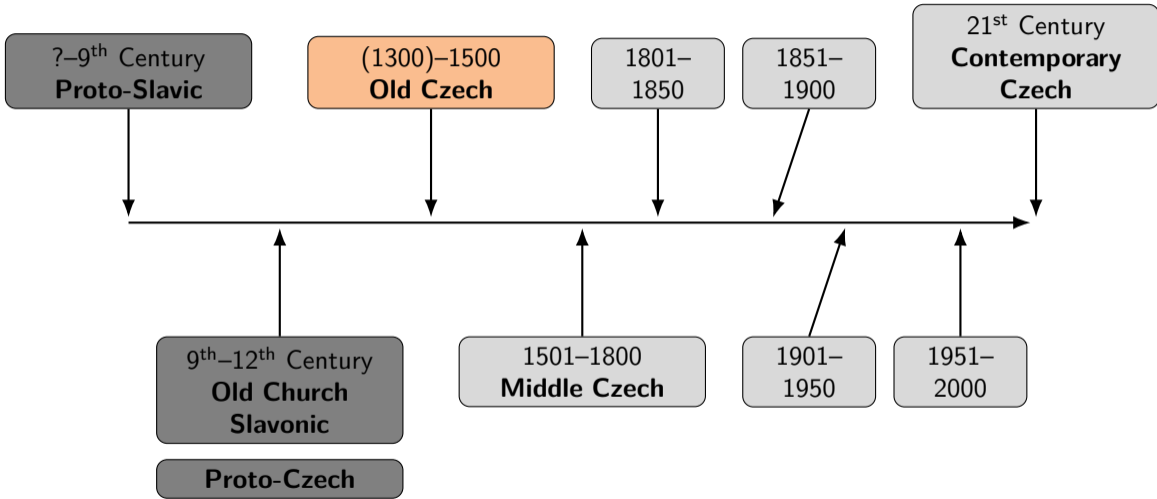
History of the Czech Language



History of the Czech Language



History of the Czech Language



Old Czech: Workflow

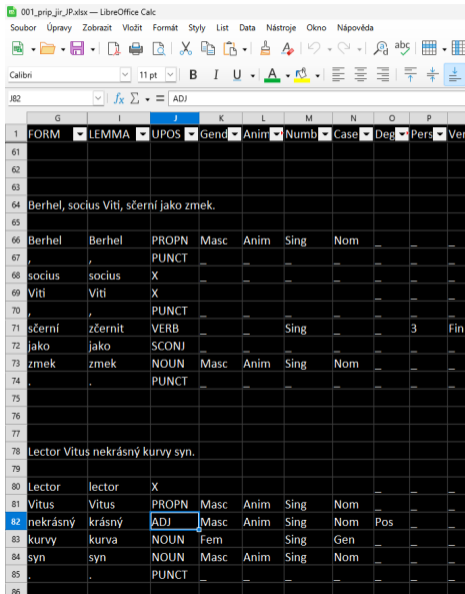
Old Czech Etalon

- 45 texts, oldest dated between 1290 and 1310
- Counts (may change during annotation):
 - 1294 sentences
 - 122,819 tokens
- Genres
 - Bible translations
 - psalms, prayers, sermons
 - epic
 - chronicles
 - regulations, law
 - medicine
 - proverbs, dream books
 - notes added to other texts

- Preprocess by **UDPipe** trained on Modern Czech (PDT)
 - Plus rule-based script that fixes some frequent words (e.g. *řka* “saying”)
 - ⇒ sentence segmentation, lemmatization, UPOS, features, dependency relations

Workflow

- Preprocess by UDPipe trained on Modern Czech (PDT)
 - Plus rule-based script that fixes some frequent words (e.g. *řka* “saying”)
 - ⇒ sentence segmentation, lemmatization, UPOS, features, dependency relations
- Convert to TSV (table)
- Annotate in Excel



001_prip_jir_JP.xlsx — LibreOffice Calc

Soubor Úpravy Zobrazit Vložit Formát Styly List Data Nástroje Okno nápověda

Calibri 11 pt B I U A

J82

	G	I	J	K	L	M	N	O	P	Ver
1	FORM	LEMMA	UPOS	Gend	Anim	Numb	Case	Deg	Pers	
61										
62										
63										
64	Berhel, socius Viti, sčerní jako zmek.									
65										
66	Berhel	Berhel	PROPN	Masc	Anim	Sing	Nom			
67	,		PUNCT							
68	socius	socius	X							
69	Viti	Viti	X							
70	,		PUNCT							
71	sčerní	zčernit	VERB			Sing			3	Fin
72	jako	jako	SCONJ							
73	zmek	zmek	NOUN	Masc	Anim	Sing	Nom			
74	.		PUNCT							
75										
76										
77										
78	Lector Vítus nekrásný kurvy syn.									
79										
80	Lector	lector	X							
81	Vítus	Vítus	PROPN	Masc	Anim	Sing	Nom			
82	nekrásný	krásný	ADJ	Masc	Anim	Sing	Nom	Pos		
83	kurvy	kurva	NOUN	Fem		Sing	Gen			
84	syn	syn	NOUN	Masc	Anim	Sing	Nom			
85	.		PUNCT							
86										

- Preprocess by UDPipe trained on Modern Czech (PDT)
 - Plus rule-based script that fixes some frequent words (e.g. *řka* “saying”)
 - ⇒ sentence segmentation, lemmatization, UPOS, features, dependency relations
- Convert to TSV (table)
- Annotate in Excel
- Two annotators per file
- Convert to CoNLL-U, postprocess, **diff**
 - ⇒ one annotator resolves differences

```
003_slx_h_AM_JP.diff.txt - Notepad2
File Edit View Settings ?
File Edit View Settings ?
1 Line 37 (Sta): Difference in UPOS: AM=AUX JP=VERB
2 Line 50 (Lydu): Difference in NameType: AM=Giv JP=Geo
3 Line 56 (Tyrie): Difference in LEMMA: AM=Tyri JP=Tyrie
4 Line 63 (jich): Difference in Gender: AM=_ JP=Masc
5 Line 63 (jich): Difference in Animacy: AM=_ JP=Anim
6 Line 64 (čakajě): Difference in Tense: AM=Past JP=Pres
7 Line 65 (ležieše): Difference in LEMMA: AM=ležiešet JP=ležet
8 Line 67 (přěčakav): Difference in Tense: AM=Pres JP=Past
9 Line 67 (přěčakav): Difference in Aspect: AM=Imp JP=Perf
10 Line 68 (tu): Difference in LEMMA: AM=ten JP=tu
11 Line 68 (tu): Difference in UPOS: AM=DET JP=ADV
12 Line 68 (tu): Difference in Gender: AM=Fem JP=_
13 Line 68 (tu): Difference in Number: AM=Sing JP=_
14 Line 68 (tu): Difference in Case: AM=Acc JP=_
15 Line 77 (minu): Difference in Variant: AM=Short JP=Long
16 Line 85 (vojšče): Difference in LEMMA: AM=vojště JP=vojska
17 Line 85 (vojšče): Difference in Gender: AM=Neut JP=Fem
18 Line 85 (vojšče): Difference in Case: AM=Acc JP=Dat
19 Line 86 (odtad): Difference in PronType: AM=Rel JP=Dem
20 Line 94 (Divno): Difference in LEMMA: AM=divno JP=divný
21 Line 94 (Divno): Difference in UPOS: AM=ADV JP=ADJ
22 Line 94 (Divno): Difference in Gender: AM=_ JP=Neut
23 Line 94 (Divno): Difference in Number: AM=_ JP=Sing
24 Line 94 (Divno): Difference in Case: AM=_ JP=Nom
25 Line 94 (Divno): Difference in Variant: AM=_ JP=Short
26 Line 100 (sdieti): Difference in LEMMA: AM=sdít JP=zdít
27 Line 102 (hna): Difference in LEMMA: AM=hna JP=hnát
28 Line 102 (hna): Difference in Aspect: AM=Act JP=Imp
29 Line 102 (hna): Difference in Voice: AM=_ JP=Act
30 Line 105 (ml): Difference in Number: AM=Sing JP=Plur
31 Line 106 (třidcěti): Difference in Case: AM=Gen JP=Acc
32 Line 122 (vzvěda): Difference in LEMMA: AM=vzvěda JP=zvědět
33 Line 129 (přijěda): Difference in LEMMA: AM=přijěda JP=přijet
34 Line 136 (tu): Difference in LEMMA: AM=tady JP=tu
35 Line 138 (Eufrates): Difference in Animacy: AM=Anim JP=Inan
36 Line 138 (Eufrates): Difference in NameType: AM=Giv JP=Geo
```

Workflow

- Preprocess by UDPipe trained on Modern Czech (PDT)
 - Plus rule-based script that fixes some frequent words (e.g. *řka* “saying”)
 - ⇒ sentence segmentation, lemmatization, UPOS, features, dependency relations
- Convert to TSV (table)
- Annotate in Excel
- Two annotators per file
- Convert to CoNLL-U, postprocess, diff
 - ⇒ one annotator resolves differences
- Official UD **validation**
- Additional consistency checks in **Udapi**
 - required & allowed features
- Return, fix, repeat...

```
# sent_id = 13_19_stol-003_alx_h-p1-s12
# text = Tak se by vypravil hrdě, vsů věcú silně i tvrdě.
Tak tak ADV _ PronType=Dem dep _
se se PRON _ Case=Acc|PronType=Prs|Reflex=Yes obj
by být AUX _ Aspect=Imp|Mood=Ind|Person=3|Polarity=Pr
vypravil vypravit VERB _ Animacy=Anim|Aspect=Perf|Gend
hrdě hrdě ADV _ Degree=Pos|Polarity=Pos dep SpaceAft
, , PUNCT _ punct _
vsů všechen DET _ Case=Ins|Gender=Fem|Number=Sing|Pro
věcú věc NOUN _ Case=Ins|Gender=Fem|Number=Sing dep
silně silně ADV _ Degree=Pos|Polarity=Pos dep _
i i CCONJ _ cc _
tvrdě tvrdě ADV _ Degree=Pos|Polarity=Pos conj Sp
. . PUNCT _ punct _
```

```
# sent_id = 13_19_stol-003_alx_h-p1-s15
# text = S obů stranů toho voza, jež jich viera i jich hróza
S s ADP _ AdpType=Prep|Case=Gen case _
obů obě NUM _ Case=Gen|Gender=Fem|Number=Dual|NumFr
stranů strana NOUN _ Case=Gen|Gender=Fem|Number=Dual
toho ten DET _ Animacy=Inan|Case=Gen|Gender=Masc
voza vůz NOUN _ Animacy=Inan|Case=Gen|Gender=Masc|F
, , PUNCT _ punct _
jež jenž PRON _ Case=Nom|Gender=Fem|Number=Sing
jich on DET _ Case=Gen|Gender=Fem|Number=Sing|F
viera víra NOUN _ Case=Nom|Gender=Fem|Number:
i i CCONJ _ cc _
jich on DET _ Case=Gen|Gender=Fem|Number=P:
hróza hrůza NOUN _ Case=Nom|Gender=Fem|Number=Si
, , PUNCT _ punct _
jedieše být AUX _ Aspect=Imp|Mood=Ind|Number=Sing|Pers
```

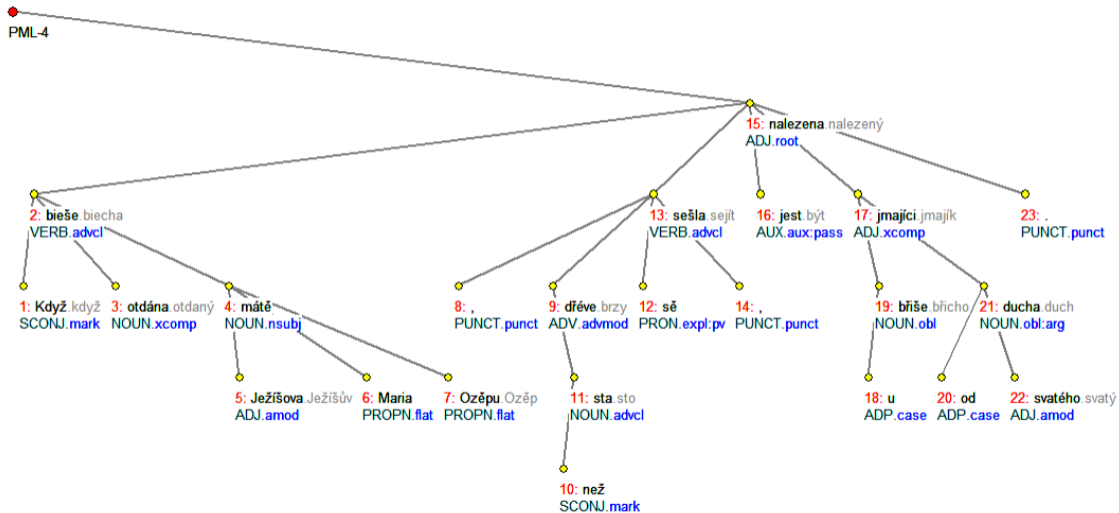
Czech UD Treebanks

- PDT (Prague Dependency Treebank)
 - Lidové noviny + Mladá Fronta + ČM Profit + Vesmír, 1993–1994
 - 87K sentences, 1.5M words
- CAC (Czech Academic Corpus / Korpus věcného stylu)
 - non-fiction, 1971–1985
 - 24K sentences, 493K words
- FicTree
 - fiction, from Czech National Corpus, 1991–2007
 - 12K sentences, 166K words
- CLTT (Czech Legal Text Treebank)
 - The Accounting Act (Zákon o účetnictví), 1991–2016
 - 1K sentences, 36K words
- PUD (Parallel Universal Dependencies)
 - online news + Wikipedia, translated from en/de/fr/it/es, around 2016
 - 1K sentences, 18K words
- Poetry
 - poetry, 19th century
 - 297 sentences, 6K words

PDT Model vs. Old Czech Data

- Genre, vocabulary: news vs. Bible
- Old vocabulary
- Orthography
 - Cleaned, transcribed, unified
 - But still not modern forms: *sě*, *viece*
- Grammar:
 - Dual number
 - Simple past (imperfect, aorist) (*bieše*, *vecě*, *jide*)
 - Converbs (přechodníky) (*řka*, *přistúpiv*)

Example Parse (UDPipe 2.0 on UD PDT 2.6)



Example Parse (UDPipe 2.0 on UD PDT 2.6)

ID	FORM	LEMMA	UPOS	XPOS	FEATS
1	Ale	ale	CCONJ	-	-
2	Kristovo	Kristův	ADJ	-	Case=Nom Gender=Neut Gender[psor]=Masc NameType=Sur Number=Sing Poss
3	porozenie	porozenie	NOUN	-	Case=Nom Gender=Neut Number=Sing Polarity=Pos
4	tak	tak	ADV	-	PronType=Dem
5	bieše	biešať	VERB	-	Aspect=Imp Mood=Ind Number=Sing Person=3 Polarity=Pos Tense=Pres Ver
6	.	.	PUNCT	-	-

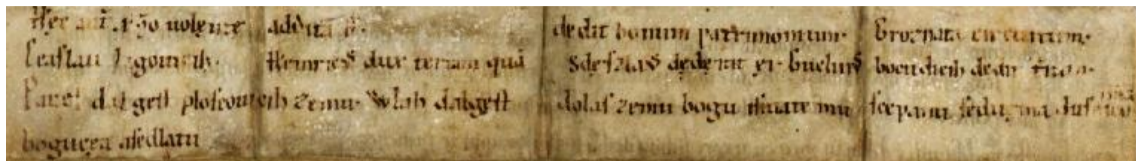
Bootstrapping

- 10 texts already annotated
 - 958 sentences, 20K words
- Retrain UDPipe
- Reparse the remaining files
- Hopefully better \Rightarrow less work for annotators
- Earlier experiments with Dresden Bible, Gospel of Matthew:

Model	Lemma	UPOS	Features
UDPipe 1.2 on PDT UD 2.5	70	77	55
UDPipe 1.2 on FicTree 2.10 + 1669 words of Matthew	78	85	65

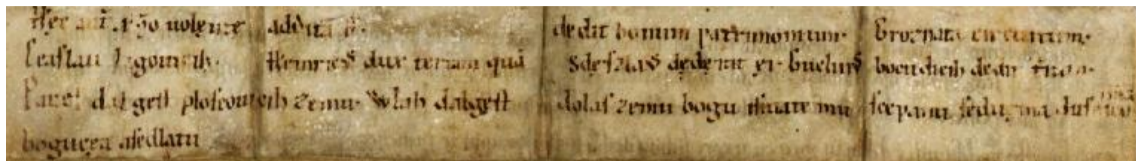
Old Czech: Interesting Issues

Modern Orthography vs. Modern Language



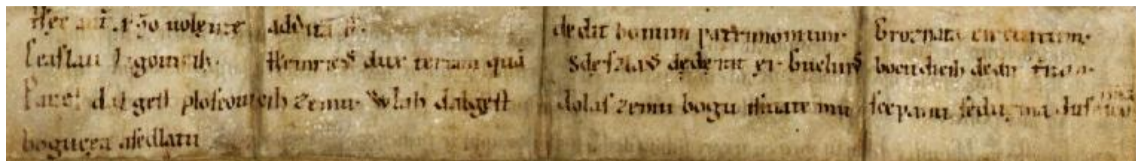
- *Pavel dal gest ploskovicih zemu Wlah dalgest dolaf zemu bogu ifuiatemu ſcepanu ſeduema duſnicoma bogucea aſedlatu*
- *Pavel dal jest Ploskovicích zem' u, Vlach dal jest Dolas zem' u bogu i sv'atému Ščepánu se dvěma dušníkoma, Bogučēja a Sedlatu.*
- *Pavel dal v Ploskovicích zemi, Vlach dal v Dolanech zemi bohu i svatému Štěpánovi se dvěma dušníky, Bogučejem a Sedlatou.*
- "Pavel gave land in Ploskovice, Vlach gave land in Dolas to God and St. Stephen with two villeins, Bogučej and Sedlata."

Modern Orthography vs. Modern Language



- Pavel dal gest ploskovicih zemu Wlah dalgest dolaf zemu bogu ifuiatemu fcepanu seduema dušnicoma bogucea afedlatu
- Pavel dal *jest Ploskovicích zem'u*, *Vlach* dal *jest Dolas zem'u* bogu *i sv'atému Ščepánu* *se dvěma dušníkoma*, *Bogučeja* a *Sedlatu*.
- Pavel dal v Ploskovicích zemi, Vlach dal v Dolanech zemi bohu i svatému Štěpánovi se dvěma dušníky, Bogučejem a Sedlatou.
- “Pavel gave land in Ploskovice, Vlach gave land in Dolas to God and St. Stephen with two villeins, Bogučej and Sedlata.”

Modern Orthography vs. Modern Language



- Pavel dal gest ploskovicih zemu Wlah dalgest dolaf zemu bogu ifuiatemu ſcepanu ſeduema duſnicoma bogucea aſedlatu
- Pavel dal **jest** Ploskovicích zem' u, Vlach dal **jest** Dolas zem' u bogu i sv'atému Ščepánu se dvěma dušníkoma, Bogučēja a Sedlatu.
- Pavel dal **v** Ploskovicích zemi, Vlach dal **v** Dolanech zemi bohu i svatému Štěpánovi se dvěma dušníky, Bogučejem a Sedlatou.
- “Pavel gave land in Ploskovice, Vlach gave land in Dolas to God and St. Stephen with two villeins, Bogučej and Sedlata.”

Tokenization

- UD concept of **multiword tokens**
 - one orthographic word → multiple morphosyntactic words
- Modern Czech examples:
 - *kdybychom* “if we would” = *když*/**SCONJ** + *bychom*/**AUX**
 - *udělalas* “you.Fem.Sing did” = *udělala*/**VERB** + *jsi*/**AUX**
 - *ses* = *jsi*/**AUX** + *se*/**PRON**

Multiword Tokens in Old Czech

Preposition + accusative pronoun:

nač = *na* *co* PronType=Rel

naň = *na* *něj* PronType=Prs

naňž = *na* *nějž* PronType=Rel

oč = *o* *co* PronType=Rel

oň = *o* *něj* PronType=Prs

oňž = *o* *nějž* PronType=Rel

...

předeň = *před* *něj* PronType=Prs

skirzěňž = *skirzě* *nějž* PronType=Rel

...

zaňž = *za* *nějž* PronType=Rel

The Particle-Clitic *-ť*

Function hard to nail down: discourse connective, emphasis, sometimes even personal pronoun *ti* “you.Dat”.

mnohémuť = *mnohému* *ť*

Examples:

- *paktť, tedyť, ješčetť, ktožť, mnohémuť, vzpomenúť, neopustítť...*

Except for lexicalized cases:

- *ať, (byť), nebť, neboť, nechť, (toť), (vždyť)*

Emphatic -ž

Narrower and less fuzzy function than *-ť* – analyzed differently!

- No splitting, kept as one word
- Lemmatized without *-ž*
- New feature: **Emph=Yes**
- But not for lexicalized cases

Examples:

- *dřevniehož, rcěmež, budž, měž, přezříž...*

Except for lexicalized cases:

- *aniž, což, jakož, jehož, kdož, když, kterýž, nikdož, protož, takž, všelicož, ...*

- Opposite situation: Old spelling with space.
 - *kdo koli* “anybody”
 - prepositional phrases vs. compound adverbs

Lemmatization

- Not only picking base form in the paradigm...
- ... but also normalization among alternatives...
- ... even after modernizing orthography!
 - *Křtitel* “Baptist”
 - *Křstitel*
 - *Krstitel*

Lemmatization

- Not only picking base form in the paradigm...
- ... but also normalization among alternatives...
- ... even after modernizing orthography!
 - *Křtitel* “Baptist”
 - *Křstitel*
 - *Krstitel*
- Modern lemma vs. old lemma:
 - forms = *otsúdí* / *otsúdie* “she / they will condemn”
 - lemma candidates = *odsoudit, odsouditi, odsúditi, odsúziti, otsoudici, otsoudit, otsouditi, otsúdici, otsúdit, otsúditi, vodsoudit, vodsouditi, vodsúditi, votsoudici, votsoudit, votsouditi, votsúdici, votsúdit, votsúditi, ...*
 - (modern) lemma = *odsoudit*

Modern Lemma Not Always Easy to Identify

- *jmě* “name” → *jméno*
- *ješutenstvie* “conceit” → *ješitenství*
- *podvstati* → *podvstat?*
- *uvrci* “cast” → *uvrhnout* ... **change of inflection class**
- *přijieti* “accept” → *přijmout*
- *cizý* “foreign” → *cizí*

- Homonyms with similar meaning:
 - *břich* vs. *břicho* “belly”
 - *obecný* “common” vs. *obecní* “of community”

Negation

Modern Czech data (PDT) **until UD 2.14**:

- **VERB** productive: *dělat* “to do” → *nedělat* “not to do”
 - Lemma = affirmative form, features: **Polarity=Neg**

Negation

Modern Czech data (PDT) **until UD 2.14**:

- **VERB** productive: *dělat* “to do” → *nedělat* “not to do”
 - Lemma = affirmative form, features: **Polarity=Neg**
- **ADJ** productive (but competing with antonyms):
hezký “pretty” → *nehezký* “unpretty, ugly”
 - Lemma = affirmative form, features: **Polarity=Neg**

Negation

Modern Czech data (PDT) **until UD 2.14**:

- **VERB** productive: *dělat* “to do” → *nedělat* “not to do”
 - Lemma = affirmative form, features: **Polarity=Neg**
- **ADJ** productive (but competing with antonyms):
hezký “pretty” → *nehezký* “unpretty, ugly”
 - Lemma = affirmative form, features: **Polarity=Neg**
- **ADV**, **DET** somewhat productive:
nedávno “recently”, *nezávisle* “independently”; *neméně* “no less”

Negation

Modern Czech data (PDT) **until UD 2.14:**

- **VERB** productive: *dělat* “to do” → *nedělat* “not to do”
 - Lemma = affirmative form, features: **Polarity=Neg**
- **ADJ** productive (but competing with antonyms):
hezký “pretty” → *nehezký* “unpretty, ugly”
 - Lemma = affirmative form, features: **Polarity=Neg**
- **ADV**, **DET** somewhat productive:
nedávno “recently”, *nezávisle* “independently”; *neméně* “no less”
- **NOUN** theoretically all are negatable (even **PROPN**!) but it is rare
 - Meaning is often shifted
 - Annotation not completely consistent
 - **Inflection:** *nezávislost* “independence”, *neschopnost* “incompetence”, *nemožnost* “impossibility”, *nedodržení* “nonobservance”, *neznalost* “unawareness”
 - **Derivation:** *nebezpečí* “danger”, *nedostatek* “shortage”, *nezaměstnanost* “unemployment”, *nemovitost* “real estate”, *nesmysl* “nonsense”, *neštěstí* “misfortune, accident”

Old Czech data and Modern Czech data (PDT) since UD 2.15:

- **NOUN** does not have **Polarity**, lemma with *ne-* \Rightarrow negation is derivation
- **All** non-pronominal **ADV** have **Polarity** (*nevelmi* “not much” – this can be also quantifier, i.e., **DET**)
- Some pronominal **ADV** have it, too (*nevždy* “not always”)

Morphological Features

Simple Past Tense: Imperfect

- Old: *Ale Kristovo porozenie tak **bieše**.*
 - **VERB** VerbForm=Fin Mood=Ind Tense=Imp Aspect=Imp Voice=Act Number=Sing Person=3 Polarity=Pos
- Modern: *S narozením Ježíše Krista to **bylo** takto:*
 - **VERB** VerbForm=Part Tense=Past Aspect=Imp Voice=Act Number=Sing Gender=Neut Polarity=Pos
- English: “Now the birth of Jesus Christ took place in this way.”

Morphological Features

Simple Past Tense: Imperfect

- Old: *Ale Kristovo porozenie tak **bieše**.*
 - **VERB** VerbForm=Fin Mood=Ind Tense=Imp Aspect=Imp Voice=Act Number=Sing Person=3 Polarity=Pos
- Modern: *S narozením Ježíše Krista to **bylo** takto:*
 - **VERB** VerbForm=Part Tense=Past Aspect=Imp Voice=Act Number=Sing Gender=Neut Polarity=Pos
- English: “Now the birth of Jesus Christ took place in this way.”

Simple Past Tense: Aorist

- Old: *Tehdy oni **pověděchu** jemu:*
 - **VERB** VerbForm=Fin Mood=Ind Tense=Past Aspect=Perf Voice=Act Number=Plur Person=3 Polarity=Pos
Variant=Long
- Modern: *Oni mu **řekli**:*
 - **VERB** VerbForm=Part Tense=Past Aspect=Perf Voice=Act Number=Plur Gender=Masc Animacy=Anim
Polarity=Pos
- English: “They told him,”

Conditional vs. Aorist of “to be”

- | | | Sing | Dual | Plur | |
|------------------------------------|---|-------------|----------------|---------------|-----------------------|
| • Aorist of <i>být</i> “be” | 1 | <i>bych</i> | <i>bychově</i> | <i>bychom</i> | *) and other variants |
| | 2 | <i>by</i> | <i>bysta</i> | <i>byste</i> | |
| | 3 | <i>by</i> | <i>bysta</i> | <i>bychu</i> | |
-
- | | | Sing | Plur |
|---------------------------------------|---|-------------|---------------|
| • Modern Conditional auxiliary | 1 | <i>bych</i> | <i>bychom</i> |
| | 2 | <i>bys</i> | <i>byste</i> |
| | 3 | <i>by</i> | <i>by</i> |
- Even modern *by* can be 2nd person in *by ses*
 - In Old Czech, aorist serves as conditional auxiliary...
 - ... but *by* is used in various persons and numbers
 - Aorist usages must be distinguished from conditional

Dual Number

- Old: ... *uzřě dva bratry, ... že biešta rybářě.*
 - **NOUN** Gender=Masc Animacy=Anim Number=Dual Case=Nom
- Modern: ... *uviděl dva bratry, ... byli totiž rybáři.*
 - **NOUN** Gender=Masc Animacy=Anim Number=Plur Case=Nom
- English: "... he saw two brothers, ... for they were fishermen."

Dual Number

- Old: ... *uzřě dva bratry, ... že biešta rybářě.*
 - **NOUN** Gender=Masc Animacy=Anim Number=Dual Case=Nom
- Modern: ... *uviděl dva bratry, ... byli totiž rybáři.*
 - **NOUN** Gender=Masc Animacy=Anim Number=Plur Case=Nom
- English: "... he saw two brothers, ... for they were fishermen."

Animacy

- Not significant grammatically as in modern Czech
- But tentatively annotated anyway, to be consistent with modern Czech data

Morphological Features

Converbs (= Gerunds = Transgressives)

- Old: *Tehdy ona ihned **ostavše** sieti, jidesta po něm.*
 - **VERB** VerbForm=Conv Tense=Past Aspect=Perf Voice=Act Number=Dual Polarity=Pos
- Modern: *Oni hned **opustili** sítě a následovali ho.*
 - **VERB** VerbForm=Part Tense=Past Aspect=Perf Voice=Act Number=Plur Gender=Masc Animacy=Anim
Polarity=Pos
- English: “Immediately they left their nets and followed him.”

Morphological Features

Converbs (= Gerunds = Transgressives)

- Old: *Tehdy ona ihned **ostavše** sieti, jidesta po něm.*
 - **VERB** VerbForm=Conv Tense=Past Aspect=Perf Voice=Act Number=Dual Polarity=Pos
- Modern: *Oni hned **opustili** sítě a následovali ho.*
 - **VERB** VerbForm=Part Tense=Past Aspect=Perf Voice=Act Number=Plur Gender=Masc Animacy=Anim Polarity=Pos
- English: “Immediately they left their nets and followed him.”

Accusative Converbs?

- Old: *... někteří ... neuzříe syna člověčieho, **přijdúce** v svém království.*
 - **VERB** VerbForm=Conv Tense=Pres Aspect=Perf Voice=Act Number=Sing Gender=Masc Animacy=Anim Polarity=Pos Case=Acc
- Modern: *... někteří ... nespatri Syna člověka **přicházejícího** ve své královské moci.*
 - **ADJ** VerbForm=Part Tense=Pres Aspect=Imp Voice=Act Number=Sing Gender=Masc Animacy=Anim Polarity=Pos Case=Acc
- English: “... some ... (will not) see the Son of Man coming in his kingdom.”
- Decision: No **Case** feature with converbs.

Supine vs. Infinitive (???)

- Old: *Nalezeny sú oslice, jíchžtos byl šel hledat.*
 - **VERB** VerbForm=Sup Aspect=Imp Polarity=Pos
- Modern: *Oslice, které jsi šel hledat, se našly.*
 - **VERB** VerbForm=Inf Aspect=Imp Polarity=Pos
- English: “The donkeys that you went to seek are found.”
- UD knows supine in Modern Slovenian.
- But what criteria would we use in Czech?

Summary

- 100K UD-style etalons to train UDPipe for each period
- Current work:
 - The ÚJČ team + DZ: Old Czech
 - The ÚČNK team + BŠ: 19th century
 - The NK team: preparing data from 20th century
- Modern lemma to facilitate searching
- Morphological features for extinct phenomena
- Tokenization: Multi-word tokens

<https://korpus.cz/hickok>