

TVORBA KORPUSŮ MLUVENÉHO JAZYKA

MILOSLAV VONDRÁČEK

Slezská univerzita v Opavě, Filozoficko-přírodovědecká fakulta v Opavě,
Česká republika
miloslav.vondracek@fpf.slu.cz

Má někdejší univerzita se v minulých letech podílela na tvorbě korpusů mluvené komunikace. Spolu se studenty jsme pořídili zvukový záznam 220 soukromých dialogických situací a jejich přepis. Při té příležitosti jsme museli řešit řadu praktických problémů. Ty vedly k formulaci podstatných teoretických otázek. K základním patří relativita jednotek mluveného jazyka. Výsledkem je metodologie tvorby korpusu, od té doby neustále zdokonalovaná. Příspěvek přináší přehled základních otázek a snaží se poskytnout teoretické odpovědi i metodiku řešení.

DOI
[https://doi.org/
10.18690/um.ff.4.2024.11](https://doi.org/10.18690/um.ff.4.2024.11)

ISBN
978-961-286-882-6

Ključne besede:
korpus mluveného jazyka,
relativita jednotek řeči,
pravidla přepisu,
zvukový záznam,
neoficiální komunikační
situace



Univerzitetna založba
Univerze v Mariboru

DOI

[https://doi.org/
10.18690/um.ff.4.2024.11](https://doi.org/10.18690/um.ff.4.2024.11)

ISBN

978-961-286-882-6

Keywords:

spoken language corpora,
relativity of speech units,
transcription rules,
audio recording,
unofficial communication
situations

CREATION OF SPOKEN LANGUAGE CORPORA

MILOSLAV VONDRÁČEK

University of Ljubljana, Faculty of Arts, Ljubljana, Slovenia
mojca.smolej@ff.uni-lj.si

In recent years, my former university has been involved in the creation of corpora of spoken communication. The students made audio recordings of 220 private conversations. With my help, the students converted these dialogues into written text. On this occasion, we had to solve some practical problems. These difficulties led to the formulation of substantial theoretical questions. The relativity of the units of spoken language is one of the fundamental ones. The result is a corpus building methodology that has been continuously improved since then. The paper provides an overview of the fundamental questions and attempts to provide theoretical answers and a solution methodology.



TVORBA GOVORJENEGA KORPUSA

MILOSLAV VONDRÁČEK

Šlezjska univerza v Opavi, Filozofska in naravoslovna fakulteta v Opavi,
Republika Češka

miloslav.vondracek@fpf.slu.cz

V preteklih letih je moja nekdanja univerza sodelovala pri ustvarjanju korpusov govornjene komunikacije. Študenti so posneli 220 zasebnih pogovorov in jih z mojo pomočjo transkribirali. Ob tej priložnosti smo morali rešiti vrsto praktičnih težav. Te težave so privedle do oblikovanja pomembnih teoretičnih vprašanj. Eno temeljnih je relativnost enot govornjenega jezika. Rezultat tega je metodologija gradnje korpusa, ki se od takrat nenehno izboljšuje. Prispevek podaja pregled temeljnih vprašanj ter poskuša podati teoretične odgovore in metodologijo reševanja.

DOI

[https://doi.org/
10.18690/um.ff.4.2024.11](https://doi.org/10.18690/um.ff.4.2024.11)

ISBN

978-961-286-882-6

Ključne besede:

korpusi govornjenega jezika,
relativnost govornnih enot,
pravila transkripcije,
zvočno snemanje,
neuradne komunikacijske
situacije



Univerzitetna založba
Univerze v Mariboru

1 Úvod

Tvůrci korpusů mluvené komunikace řady ORAL, vznikajícího úsilím Ústavu Českého národního korpusu Filozofické fakulty Univerzity Karlovy v Praze, se zhruba v roce 2005 obrátili na další bohemistická univerzitní pracoviště České republiky s výzvou k účasti na jejich projektu.¹ Pracoviště už tou dobou mělo zkušenost s tvorbou Pražského mluveného korpusu (2001, dále PMK²), byly tu poznatky spjaté s Brněnským mluveným korpusem (2002, dále BMK³), z aktuálního sběru materiálu pro celočeský ORAL2006,⁴ to vše na pozadí důkladné obeznámenosti s problematikou v mezinárodním měřítku. Spolupráce s mimopražskými středisky měla spočívat ve shromáždění zvukových záznamů neoficiálních komunikačních situací dostatečného rozsahu a jejich přepisu podle stanovených pravidel. Povaha spolupráce s organizačním centrem se mohla případ od případu lišit v závislosti na odborném (a nakonec i lidském) profilu prostředníka. Zpětně mohu konstatovat, že pro mne – a průkazně i pro angažované studenty⁵ – se stala tato aktivita zdrojem širokého spektra poznatků ze všech lingvistických disciplín. To je důvod, proč neváhám zkušenost s odstupem a nadhledem uplynulých let znovu sdílet.⁶

2 Povaha komunikační situace

Korpusy ORAL2008 a ORAL2013,⁷ k nimž se váže má akviziční zkušenost, jsou charakterizovány jako korpusy neformální mluvené češtiny. Sdružují „materiál představující prototypický spontánní mluvený jazyk, který se používá při bezprostřední interakci mluvčích v neformálních komunikačních situacích.

¹ Výsledky se promítají do korpusů ORAL2008 a ORAL 2013. K nim více na <https://wiki.korpus.cz/doku.php/cnk:uvod>

² Informace o PMK dostupné na <https://wiki.korpus.cz/doku.php/cnk:pmk>

³ Informace o BMK dostupné na <https://wiki.korpus.cz/doku.php/cnk:bmkm>. PMK i BMK ovšem zahrnují zčásti i moderované formální monologické promluvy.

⁴ Informace o ORAL2006 dostupné na <https://wiki.korpus.cz/doku.php/cnk:oral2006>.

⁵ Mé tehdejší působiště: Univerzita Hradec Králové, Pedagogická fakulta, roky 2005–2011. Zapojilo se postupně přes 200 studentů. Každý pořídil jedinečnou sondu do mluvené komunikace, tj. nahrávku o délce průměrně 30 minut (reálně 20–50 minut, celkem přes 100 hodin). Student sám svůj zvukový záznam přepsal (v první etapě do textového editoru Word, později do aplikace Transcriber). Přepis poté prošel několikanásobnou kontrolou věrnosti přepisu. V prvních kolech kontrolu zajišťoval příslušný akademický pracovník univerzity (na UHK já), finální kontrolu pracovník ÚČNK FF UK. U obou spolutvořených korpusů viz úsek Poděkování (<https://wiki.korpus.cz/doku.php/cnk:oral2008>, <https://wiki.korpus.cz/doku.php/cnk:oral2013>).

⁶ O aktuálních zkušenostech z probíhajícího projektu jsem hovořil na konferenci Čestina v mluveném korpusu, Praha 2007, pro více informací viz Vondráček 2008.

⁷ Oba ve srovnání s korpusem ORAL2006 kromě Čech zahrnují i oblast Moravy a Slezska (což bylo jedním z motivů k oslovení mimopražských univerzitních pracovišť).

Hlavními kritérii pro získávání nahrávek byly: fyzická přítomnost všech mluvčích na jednom místě, dialogičnost promluv, vzájemný blízký vztah mluvčích, nepřipravenost, spontánnost, neveřejná a neoficiální komunikační situace“ (ORAL2013, kráceno).

Starší korpusy popisují neformální komunikační situace dalšími znaky: „neformální promluvy tvoří dialogy dvou, případně i více mluvčích, kteří se dobře znají“ (BMK), resp. „kteří se znají“ (PMK). V případě korpusu ORAL2006 „[v]šechny nahrávky vznikaly výhradně v neformálních situacích, což znamená, že se mluvčí vzájemně znali a měli k sobě přátelský vztah.“ Při tomto zadání se projevila nejednoznačnost v interpretaci vykání (i mezi rodinnými příslušníky). Současně se objevil jiný nečekaný situační rys – účastník komunikace je sice v přátelském vztahu ke komunikačnímu partnerovi, současně však hovoří z titulu své profese (kamarádka kadeřnice při úpravě účesu, kamarád automechanik při sjednávání opravy vozu). Průnikem obou eventualit vzniklo další pomocné kritérium, v dostupných podkladech nezaznamenané: vykání samo o sobě není v rozporu s neformálností a neoficiálností komunikační situace; současně (bez ohledu na tykání) žádný z účastníků nevystupuje v dialogu z titulu své profese. Tato charakteristika se více než osvědčila při jemném rozlišení komunikačních situací v rámci stylistických analýz komunikační sféry běžné.

3 Delimitace syntaktických jednotek

Pro první přepisy záznamů (pořizovaných již na digitální záznamníky) jsme využívali textového editoru Word.⁸ Ještě pro předchozí korpus ORAL2006 platila pravidla stanovení hranic syntaktických celků víceméně odpovídající zvyklostem psaných textů.⁹ Mírný posun přináší korpus ORAL2008 ve způsobu zaznamenání obvyklých defektů mluveného projevu.¹⁰

⁸ Složitou cestu uvědomování si a dokumentace rozdílů psané a mluvené češtiny ilustruje text J. Hoffmannové a M. Míkulecké (2011, 78–92).

⁹ „Hranice vět se vyznačují jen interpunkcí, na začátku věty píšeme malé písmeno. Větné interpunkce se užívá tak, jak je to obvyklé v textech psaných, tj. nezachycuje se přerušování věty pauzami, naopak náležitou čárku ve větách a v souvětích píšeme, i když se větné předěly pauzou nerealizují. Neukončené věty označujeme třemi tečkami s dvojtečkou odsazenými mezerou, tedy ...: Příklad, kdy je mluvčí přerušen jiným mluvčím, ale ve výpovědi později pokračuje, se v dialogu značí třemi tečkami na konci přerušené výpovědi i na začátku její navazující části.“ (archiv autora, viz též dále) „Pořizování nahrávek, jejich přepisování a označování probíhalo v souladu s obecnými zásadami uplatňovanými při přípravě všech předchozích mluvených korpusů v rámci Českého národního korpusu, zejména korpusu ORAL2006.“ (<https://wiki.korpus.cz/doku.php/cnk:oral2008>)

¹⁰ „Zvláštní rysy syntaktické stránky mluvených projevů (přeřeknutí, přerušování a změny větné perspektivy, přiřazování vět a větných úseků apod.) jsou zachyceny zjednodušeně, většinou pouze pomocí čárky. Pokud se slova opakují, jsou oddělena čárkou.“ (ORAL2008)

<3> *dovopravdy funguje.*

<2> *já vim, nó, my sme jak já, já, né my, co budu říkat my, já sem za vola.*

<3> *né, to není, není, já ti říkám, čím to bylo, vy ste podle mě tady s tím šabo, šibovali a naklonilo se to, udělala se bublina a to je jedna z alternativ, která byla možná, a pře, pak to nechladilo, protože si to musí sednout.*

<2> *ale vono, já mam takovej pocit, že s tím předtím jakoby nikdo nebejbal a najednou to přestalo jít.*

<3> *druhá varianta je, že, že se, žes v podstatě se to dostalo do nějaký mezipolohy a že, eee, třeba při tom přepínači, což mohlo být, a že začal blbnout termostat. a znovu začal fungovat. z nákejch důvodů.*

<2> *hm, hm.*

<0> *to je nejspíš:*

<2> *to je, to sou, todle, Jarečku, viš, to mi vysvětlí, viš, že já se tady ztrapňuju, rok tady máme ledničku vypnutou ... (úryvek ze sondy 05H006N)*

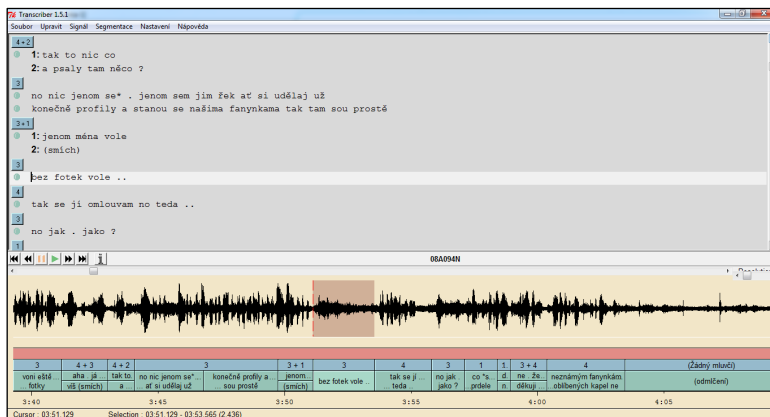
Počínaje korpusem ORAL2013 se přepis odehrává v transkripčním programu Transcriber.¹¹ To s sebou přináší – i vzhledem ke změně techniky převodu, opoře ve zrakové kontrole zvukového záznamu, odlišné segmentaci toku řeči¹² – změnu vnímání toku dialogu:¹³ systematictější a komunikačně velmi přínosně se sledují překryvy komunikantů¹⁴ (opouští se tedy snaha rozplést je do zdánlivě samostatných, nesimultánních replik), a zejména interpunkce převzatá z psaných textů je nahrazena interpunkcí pauzovou. Odlišují se tři kvality pauzy: pauza krátká, pauza delší a odmlčení; rozlišený jsou podle individuálního tempa jednotlivých mluvčích.

¹¹ Program Transcriber neslouží k automatickému převodu zvukového záznamu na psaný text. Po uložení nahrávky ale přepisovatelé umožňuje členit text do krátkých sekvencí podle kritérií formálních (max. délka, viz dále) nebo funkčních (replika jednoho mluvčího, pasáž překrývání více mluvčích ap.) Program sloužil jen pro přepis a kontrolu. Hotová, zkontrolovaná sonda převzatá pro zařazení do korpusu se dále zobrazuje v jiném nástroji (KonText, <https://www.korpus.cz/kontext/query?corpname=syn2020>, viz obr. 2) Všechny sondy vytvořené na Univerzitě Hradec Králové pod mou patronací mají na třetí pozici sedmimístné značky sondy písmeno H (př.: 09H015N = sonda vytvořená v roce 2009 v Hradci Králové jako 15. v pořadí, komunikační situace neoficiální).

¹² „[r]epliky jsou členěny na segmenty představující sémanticky, prozodicky i syntakticky ucelenou sekvenci v průměru o 5–10 slovech (maximálně však 15).“ (ORAL2013)

¹³ K vertikální ose komunikátu viz kupř. podnětnou studii Nepravá hypotaxe v spontánních mluvených projevech (Bílková 2021)

¹⁴ K tomu např. Komrsková–Poukarová 2018, 41–56; Komrsková–Poukarová–Havlík 2019, 102–116.



Obrázek 1: Přepis zvukového záznamu v programu Transcriber

Zdroj: <https://wiki.korpus.cz/doku.php/cnk:oral2013>. Printscreensy aplikace s přepisem studentských sond jsem nepoživoval

Tradiční interpunkcí se vyznačuje pouze nápadně stoupavá intonace věty tázací, popř. vykřičníkem nápadná zvolací intonace. Důvodem rezignace na standardní značení mezivětných předělů bylo zejm. stanovisko odborníků z oblasti fonetiky a fonologie,¹⁵ že zvukové hranice takových vyšších textových jednotek lze nalézt jen stěží. Opakovaně se nám tak potvrzuje, že sémaziologická dekompozice textu probíhá primárně na základě kritérií významových.

Tabulka 1: Ukázka přepisu dialogu s pauzovou interpunkcí

Radomíra_7878	–	<i>to úplně cejtím že mně to chybí ty jo su úplně taková ..</i>
Adéla_5592	–	<i>no já taky . chcu jako něco . pořádně .</i>
Adéla_5592	–	<i>když -</i>
+ Radomíra_7878	–	<i>mně hlavně .</i>
Radomíra_7878	–	<i>já hlavně cejtím jak sem závislá na Dimim . až teď s* dycky když jakoby sem vod něj .</i>
Adéla_5592	–	<i>@</i>
Radomíra_7878	–	<i>teď jak sem sem jela tak sem říkala . ty jo mně se sem nechce .</i>
Adéla_5592	–	<i>(smích)</i>
Radomíra_7878	–	<i>najednou musím všechno sama vš všechno</i>
Adéla_5592	–	<i>no jasné no</i>
+ Radomíra_7878	–	<i>tak jakože ono ale</i>
Radomíra_7878	–	<i>mně to nepřide</i>

Vir: Zdroj: Korpus ORALv1, sonda 11A085N, <https://www.korpus.cz/>

¹⁵ Fonetický ústav Filozofické fakulty Univerzity Karlovy Praha, tehdy pod vedením prof. PhDr. Zdenky Palkové, CSc. (<https://fonetika.ff.cuni.cz/ustav/vyucujici/zdena-palkova/>)

4 Fenomén lingvální, paralingvální a nonlingvální

Nikoli nezajímavé je (a bylo nejen pro studenty filologických oborů) ověřování hranic žvlu lingválního a nonlingválního s přechodovou oblastí v jevech povahy paralingvální. Pravidla platná pro korpus ORAL2006 stanovila dosti jednoduše, že parazitní zvukové projevy se zaznamenávají sekvencí tří písmen: obvykle *hmm* pro zvuky spíše souhláskové, *eee* pro zvuky povahy převážně vokaliké. Smích se vyznačuje poznámkou v závorce (smích); v případě verbalizovaného, resp. i verbalizovatelného smíchu se zapisuje zvuková forma „co nejbliže slyšenému“.¹⁶

Už v průběhu sběru jazykového materiálu pro korpus ORAL2008 jsme narazili na nezbytnost odlišení komunikačně relevantních neartikulovaných nebo poloartikulovaných zvuků vyjadřujících souhlas, nesouhlas, pochybnosti, váhání apod. Zejména šlo o laryngálně-nazální zvuk (provázený případně odmítavým pohybem hlavou), graficky snad zaznamennatelný jako kombinace rázu a písmena pro konsonant *m*, tj. jako *'m-'m*. (Dlužno říci, že např. ani autoři Příruční mluvnice češtiny, dále PMC, nečinili nijak podstatného rozdílu mezi výrazy *hm*, *ehm* a *mbm*, radice je souběžně k interjekcím (první dva) a partikulím (třetí), obecně přitakacím (PMČ 1995, 356 - § 593, a 365 - § 606). Garanti nově vznikajícího korpusu pochopili důležitost signalizace zásadně odlišné komunikační funkce prostředků hraniční povahy a stanovili další, strukturovanější pravidla transkripce: *hmm* pro přitakací responzní zvuk, *emem* pro nesouhlasný responzní zvuk, *mmm* pro souhláskové hezitační zvuky a *eee* pro zvuky hezitační samohláskové.¹⁷

Nejde o jev funkčně triviální. Při tvorbě hesla *ehm* pro Akademický slovník spisovné češtiny (dále ASSČ) odlišila Jana Špirudová¹⁸ sedm jeho různých komunikačních platností (v dnešní zveřejněné verzi jsou dostupné čtyři, pročež uvádím kompletní znění hesla v neoklešněné podobě):

ehm

- V1 vyjádření nejistoty, váhání, nerozhodnosti, mírných rozpaků
Ehm, stala se taková věc...; Já jsem, ehm – zkrátka, neřlobte se, že jsem vás tak otravoval.; Ehm, jak ti to říct, zkrátka kluci nevěděli, že ta bonboniéra není naše.; Ehm, nevím, vážně mě to nijak extra neláká.; „Mohl bych... Ehm, mohl bych vás o něco požádat?“ vykotal.

¹⁶ Zásady přepisu pro korpus ORAL2006. Dostupné z: https://wiki.korpus.cz/doku.php/seznamy:pravidla_2006

¹⁷ Inspirativní je zpracování výrazu *no* jako diskurzivního markeru (Bílková-Zeman 2020, 191-198)

¹⁸ Jako spoluautor ASSČ a osoba pověřená kontrolou tohoto hesla vidím informace o zakladateli (zpracovateli) hesla a autorech jeho vyšších verzí. Protože interjekce představují jednu z mých gramatických specializací, Jana Špirudová mě o revizi tohoto hesla výslovně požádala.

- V2 opatrný, váhavý souhlas, kladná odpověď na otázku
Dáte si pivo? – Ehm, díky.; Ehm, tak to můžu říct docela přesně.
- V3 zdůraznění otázky, pobídka k odpovědi
Ehm, nevádí ti, když půjdu s vámi?; Ehm a jak byste sestavil rozpočet vy?; Moby vám, ehm, ehm, nějak pomoci?
- V4 zvuk při odkašlání, kterým chce zprav. někdo na něco upozornit (např. že se chystá něco říct)
„Ehm,“ odkašlal si ředitel.; „Ehm, ehm, ehm,“ zakašlala paní Renata a kopla bo pod stolem prudce do holeně.; „Ehm, ehm,“ odkašlala jsem si významně, jak to mám ve zvyku před klíčovým projevem na jednání s komplikovaným partnerem.
- V5 výraz zamyšlení nebo snahy posoudit návrh, event. rozvinout téma
Ehm, co kdybychom šli radši dovnitř, je tu celkem zima.; Ehm, ehm...ještě jedna věc.; Ehm – a co takhle použít polarizační filtr?
- V6 výraz nelibosti, nesouhlasu nebo mírného odmítnutí
Ehm, to teda děkuju.; Ehm, to jste myslel vážně, pane? Ehm, to jako nemám říkat?
V7 povzdechnutí, vyjádření mírné námitky
Ehm, a je to tady zase.; Ehm, tak to už tady jednou bylo.

V úzu se vedle základní podoby citoslovce *ehm* objevují varianty opakující písmeno *m*, čímž se naznačuje delší trvání zvuku a zesiluje se tak komunikační účinek: *ehmm, ehmmm*. (ASSČ, heslo *ehm*, verze z 26. 5. 2022, před umrtvením významů 5-7)

Tabulka 2: Frekvence grafických forem rezponzních a hezitačních zvuků

	Grafická forma	Absolutní frekvence ¹⁹	Frekvence na milion slov (i. p. m.)
ORALv1	<i>hmm</i>	50 095	7 874,46
	<i>emm</i>	636	99,97
	<i>ééé</i>	5	0,79
	<i>em em</i>	1	0,16
	<i>chm</i>	1	0,16
PMK	<i>hm</i>	760	927,66
	<i>ehm</i>	143	174,55
	<i>eee</i>	60	73,24
	<i>hmm</i>	35	42,72
	<i>chm</i>	27	32,96
	<i>mmm</i>	25	30,52
	<i>ééé</i>	4	4,88
	<i>emem</i>	3	3,66
	<i>em em</i>	1	1,22
	<i>ehmm</i>	1	1,22
ORTOFON	<i>hmm</i>	29 200	11 403,62

¹⁹ Absolutní frekvence je ovlivněna velikostí konkrétního korpusu. Parametr i. p. m. udává přepočtenou frekvenci na ideální jeden milion slov daného korpusu.

	Grafická forma	Absolutní frekvence ¹⁹	Frekvence na milion slov (i. p. m.)
	<i>emm</i>	579	226,12
	@ / @@	23 606	9 218,97
	<i>chm</i>	6	2,34
	<i>eee</i>	3	1,17
	<i>mmm</i>	1	0,39

Zdroj: vlastní

S ohledem na budoucí lemmatizaci byla cíleně sjednocována přirozeně oscilující forma zápisu (respektující nakonec pokyn zapisovat co nejvěrněji slyšené). Logicky se tak potlačuje možnost vyjádřit počtem grafémů či jejich skladbou délku či obecnou povahu zvuku. Do transkriptů se však přece jen individuální reflexe v jisté míře prosadila.

Současná praxe dvouúrovňového přepisu, uplatňovaná pro korpus ORTOFON, nahlíží tuto oblast komunikace zase jiným prizmatem. Pravidla²⁰ ukládají přepisovat na rovině ortografické responzní zvuky přitakávací *hmm*, responzní zvuky odmítací *emm*, citoslovce bez lexikalizované podoby jako & a hezitací zvuky znakem @ pro kratší hezitace a @@ pro delší. Není zřejmé, jestli ojedinělé odlišné zaznamenané formy vypovídají o snaze vědomě zachytit specifický jev přesahující pravidla, nebo jsou jen projevem přehlédnutí. (Pro srovnání: Korpusy psaného jazyka dokládají, že psané texty, zejm. beletristické a publicistické, pracují i s dalšími formami, v mluvených korpusech registrované nanejvýš okrajově; v SYNv11 kupř. *chm* 561x / 0,09 i.p.m., *chmm* 131x / 0,02 i.p.m., *chmmm* 4x, *chmmmmm* 1x. To může odpovídat komunikačně funkčnímu, převážně nazálnímu nelaryngálnímu výdechovému zvuku, povzdechu se sevřenými rty.)

*Víte ... **chm** ... zrovna než jste přišel ... prohlížel jsem podložní skříčka ... a přišel jsem najedno, které by vás mohlo zajímat.* (korpus SYNv11)

*. --- brácha a tatínek si dycky přál syna a syna teda měl . i když se narodil pozdějc než já . ale dycky já sem byla taková .. **mmm** .. no nechci říct úplně v koutku . ale Jirka dycky moh víc ---* (korpus ORALv1)

²⁰ Transkripcie v korpusu ORTOFON. Dostupné z: <https://wiki.korpus.cz/doku.php/cnk:ortofon:pravidla>

Tabulka 3: Ukázka přepisu dialogu s pauzovou interpunkcí a responzními zvuky

Ivana Š.	– <i>bmm</i>
Rozálie B.	– <i>víš ? .. hele fakt do mě nic nebylo hele Ilonko</i>
Ivana Š.	– <i>bmm ..</i>
Ivana Š.	– <i>bmm .. bmm .. bmm</i>
Rozálie B.	– <i>do mě nic nebylo jo</i>
Ivana Š.	– <i>bmm</i>
Rozálie B.	– <i>no a naštěstí teda když přijeli ..</i>
Rozálie B.	– <i>tak mi řekli že to byl teror že je . v devět večer přivezli na pokoje</i>
Rozálie B.	– <i>a v sedm ráno vstávali protože se zase jede ne ..</i>
Ivana Š.	– <i>jasně byli urvaný</i>

Zdroj: Korpus ORTOFONv2, sonda 18X096N, <https://www.korpus.cz/>

Tabulka 4: Ukázka přepisu s nelexikalizovanými interjekcemi (&) a hezitačními zvuky (@)

Anna V.	– <i>tady v Brně má @ Kokino se to jmenuje je to značka . to dělá nějaká slečna a ..</i>
Anna V.	– <i>má tam @ čokolortik s mořskou solí měla #s to někdy ?</i>
Miriám Š.	– <i>ehm</i>
Anna V.	– <i>jestli máš ráda čokoládku a karamel .. z tohoble by ses úplně potento protože -</i>
Miriám Š.	– <i>musíš dát potom nějaký odkaz</i>
Anna V.	– <i>&</i>
Anna V.	– <i>& . to je víc kde to mají ? . to mají na Gor* @ Gorkého . víc kde</i>

Zdroj: Korpus ORTOFONv2, sonda 14A004N, <https://www.korpus.cz/>

5 Delimitace lexikálních jednotek a jejich hláskové struktury

S relativitou textových jednotek na úrovni věty a souvětí, resp. výpovědi, jakož i s otázkou hranic jevů lingválních a nonlingválních souvisí relativita centrální jednotky lexikálního systému (a svým způsobem centrální jednotka jazyka vůbec), totiž slova. Znovu platí, že vyčlenění jednotky v proudě řeči je značně relativní, formální delimitace vůči okolním jednotkám pauzou je hypotetická. Přírozený jazyk ze své podstaty inklinuje k jednotě formy a funkce, což posiluje tendenci k formálnímu splývání původně analytických jednotek. Podobně vede změna funkcí k proměnám formy. Všechny náznaky těchto procesů v mluvené komunikaci má proto smysl registrovat.²¹

²¹ K tomu viz např. Blatná 2006, 7–19, Vondráček 2013, 79–89, Volín 2015, 44–47, Vondráček 2018, 35–110.

Tak by bylo možno sledovat ukázkovou proměnu desubstantivních forem kontaktového prostředku citoslovečného a částicového rázu *ty vole* (chápaného nejen ve školském prostředí funkčně neadekvátně jako vulgární, resp. zhrubělé²² oslovení) ve zvukové a grafické realizaci: jeho hláskovou variabilitu či erozi včetně intonačního, resp. obecně prozodického pozadí v proměnách doby a v závislosti na funkci (→ part. *tyvoe* → *tyoe* → *tye*; → interj. *tývoe* → *týe* vedle snad tabuové náhražky *týjo* → *týo*). Z dobrých důvodů ovšem tematizují samostatně dva jiné jevy podobné povahy.

Důvody pro upuštění od detailní analýzy jevu již naznačeného (deskripce slovnědruhově transpozice kontaktové fráze spočívající v pronominu + vokativu substantiva *ty vole* v kontaktovou partikuli *tye* či podivovou interjekci *týe*) jsou naznačené povahy axiologické (byť příznak vulgárnosti či zhrublosti v jazykovém materiálu nebyl v žádném případě důvodem k vyřazení sondy). Volím však dva jiné, analogické, ovšem různě složité případy funkční a slovnědruhově proměny provázené destrukcí formy: proměnu substantivního vokativu *člověče* v jednotku pronominálního typu *ček*, v kontaktovou partikuli *čěče* a v podivovou interjekci *čěče*. Dílem jiné (tabuové, ovšem v tomto případě nábožensky tabuové) jsou pohnutky k obměnám forem (a k jejich následným deformacím) u interjekcí s funkcí povzdechu či zaklení typu *proboba*, *prokerindapána*. Podobných ilustračních typově odlišných dokladů je možno zvolit desítky; uvádím právě jen některé z těch, jež si obvykle i v přepisovacích pravidlech vyžádaly samostatné pojednání. K oběmu viz např. Laubeová 2020.

5.1 Typ *člověk* / *ček*, *člověče* / *čoveče* / *čěče*

Starší české slovníky zaznamenávají výraz *člověče* jako vokativ substantiva, funkčně jej označují jako expresivní oslovení, zařazují jej k hovorové nebo obecné češtině. ASSČ uvádí původní vokativní formu jako samostatné heslo, slovnědruhově interpretované jako interjekci (kolokviální vyšší), k ní dále i variantní formy *čoveče*, *čěče* jako kolokviální. Právě o zaznamenání hláskové eroze slova při slovnědruhově transpozici nám jde (a zmínit je třeba i funkci částicovou, v níž výraz nevykazuje znaky větého ekvivalentu). ASSČ také pod substantivním heslem *člověk* uvádí uplatnění ve funkci zájmena neurčitého, záporného nebo osobního s významem 'blíže neurčená osoba, někdo, kdokoliv, každý, mluvčí (já), adresát (ty, vy) apod., ve

²² SSJČ i SSČ svorně: zhrub., i nadávka, hlupák; kontaktovou frází *ty vole* však tyto zdroje ani IJP neuvádějí.

spojení se záporným slovesem žádná osoba, nikdo!. I tato slovnědruhová transpozice je v hovoru provázána destrukcí formy (*čěke*).

The screenshot shows the KonText web interface. At the top, there is a search bar with the text 'Dotaz Korpusy Uložit Konkordance Filtr Frekvence Kolokace Zobrazení Návoděda'. Below the search bar, there are several filters and options. The main part of the image shows a list of search results, each with a checkbox and a snippet of text containing the word 'čěče'. The results are sorted by frequency, with the most frequent results at the top. The text snippets are highlighted to show the context of the word.

Obrázek 2: Ukázka konkordancí s výrazem *čěče* ve webovém rozhraní KonText

Následující tabulka ukazuje, které z forem zachytil který z korpusů, popř. jak jsou morfologicky značkovány.

Tabulka 5: Korpusově evidované formy výrazů *člověk* / *ček*, *člověče* / *čoveče* / *čěče*

	Grafická forma	Absolutní frekvence	Frekvence na milion slov (i. p. m.)
PMK	<i>čoveče</i> / <i>čověče</i>	49 / 3	59,81
PMK	<i>člověče</i>	15	18,31
PMK	<i>čoeče</i>	1	1,22
PMK	<i>ček</i>	1	1,22
ORAL2008	<i>čoveče</i> / <i>čověče</i>	130 / 14	96,33
ORAL2008	<i>člověče</i>	16	11,86
ORAL2008	<i>čěče</i>	2	1,48
ORAL2008	<i>ček</i>	2	1,48
ORAL2013	<i>čoveče</i> / <i>čověče</i>	547 / 10	166,49
ORAL2013	<i>člověče</i>	16	4,87
ORAL2013	<i>čěče</i>	17	5,17
ORAL2013	<i>ček</i>	0	0
ORTOFONv2	<i>čoveče</i> / <i>čověče</i> (subst.)	12 / 0	4,69
ORTOFONv2	<i>člověče</i> (subst.)	233	90,99
ORTOFONv2	<i>čěče</i> (subst.)	0	0
ORTOFONv2	<i>ček</i> (subst.)	0	0

Zdroj: vlastní

... --- sou takový ještě vomámený ty ryby , že jo ? ... povídal, to sem ale m* ...: --- voni sou přibblý, **čoveče**, voni vůbec nejdou, voni zůstanou u břehu v tý trávě a vůbec nikam nejdou . no (korpus ORAL2008)

už sem nemohla . no a já říkám . zas to máš blbě . **čěče** co t* jako něco takovýho no a . tam je jedna taková šileně . mmm . šilená poděška . a vona botová úplně jako . **čěče** jo a toto . a teď se tam řebtali . brozný (korpus ORAL2013)

... no prostě, stálo jim to za to, no . jesi ... : no pokavaď, že jo, pokavaď je to fyzicky a duševně vodromná, tak **čěk** nemůže říct: stálo to za to, no . holt, já nerim, no . (korpus BMK)

... v noci v noci dyž se **čěk** probudí tak to bolí, a nevím jak si lehnout, fakt je že pak si nebudu smět dát noby křížem že jo třebaš, vsedě určitě ne, jo jak s* **ček** normálně sedí a přebodí si nobu tak (korpus ORALv1)

Stejně jako koreluje změna funkce slovní formy s omezením či ztrátou paradigmatu a s erozí fonetické struktury, dochází ke změnám u víceslovných lexikálních jednotek (historicky nejspíš takto haplogicky vznikla slova *vašnosti* ← *Vaše Jasnosti* s dodatečným Nsg *vašnosta*, *slečna* ← *slechtična*, tímž směrem se ubírala forma *pančelka* ← *paní učitelka*, v mírné podobě nakonec i dosud nekodifikované *nashledanou* ← *na sbledanou*, **sbledaná*). Také toto splývání forem a redukční procesy má smysl sledovat. I pro druhé uvedu ilustrační příklad.

5.2 Typ *probaha* / *probůh*, *prokristapána*, *propánajána*, *propánaboha*

Všechny velké slovníky češtiny 20. století registrují lexikální formu *probaha* (některé i *probůh*) jako citoslovce;²³ neuvádějí varianty psané analyticky (tudíž nemusí řešit pravopis velkých písmen vázaný na tuto funkci). Lístková excerpta²⁴ tvořící podklad Příručního slovníku jazyka českého (dále PSJČ) registrují ve třinácti dokladech i lexikální formy *prokristapána* a jeho tabuovou obměnu *prokřindapána*, vedle ní s nižším počtem dokladů i *propánajána*, *propánaboha*, vše hodnoceno jako interjekce (slovníková hesla ovšem nevznikla). Pozdější slovníky už heslo uvádějí (SSČ²⁵ *prokrista(pána)* i *pro křista pána*, SSJČ²⁶ *prokřista*, *prokřistapána*, psáno též *pro křista pána*, obdobně *propána*, *propánajána*, *propánaboha*, *propánakrále* (psáno též *pro pána boha*), rovněž jako interjekci. Psaní malého písmena lze přiřknout spíše režimním pohnutkám než snaze odlišit rouhavé brání jména Syna Božího nadarmo od vroucího obvolávání se na něj; proti tomuto směru uvažování naopak hovoří připouštěné analytické psaní. Nezdá se nepodstatné, jak je týž výraz zaznamenaný v přepisech pro mluvené korpusy sto let od vzniku PSJČ.

²³ Srov. též např. Kleňhová 2010; 2012, 238–254.

²⁴ Elektronicky jsou dostupná z webu Ústavu pro jazyk český AV ČR, v. v. i.: <https://psjc.ujc.cas.cz?>

²⁵ *Slovník spisovné češtiny pro školu a veřejnost*, 1994. Praha: Academia.

²⁶ *Slovník spisovného jazyka českého*, 1989. Praha: Academia.

Tabulka 5: Korpusově evidované formy výrazů *prokřistapána / propánajána / propánaboha*

	Grafická forma	Absolutní frekvence	Frekvence na milion slov (i. p. m.)
PMK	<i>proboba / probůh</i>	7 / 0	8,54 / 0
PMK	<i>pro boba</i>	0	0
PMK	<i>prokřistapána / propánajána / propánaboba</i>	3 / 0 / 0	3,66 / 0 / 0
PMK	<i>pro křista pána / pro pána jána / pro pána boba</i>	0 / 0 / 0	0 / 0 / 0
BMK	<i>proboba / probůh</i>	9 / 1	15,1 / 1,68
BMK	<i>pro boba</i>	0	0
BMK	<i>prokřistapána / propánajána / propánaboba</i>	0 / 0 / 0	0 / 0 / 0
BMK	<i>pro křista pána / pro pána jána / pro pána boba</i>	0 / 0 / 0	0 / 0 / 0
ORAL2008	<i>proboba / probůh</i>	6 / 0	4,45 / 0
ORAL2008	<i>pro boba</i>	0	0
ORAL2008	<i>prokřistapána / propánajána / propánaboba</i>	4 / 0 / 0	2,96 / 0 / 0
ORAL2008	<i>pro křista pána / pro pána jána / pro pána boba</i>	0 / 0 / 0	0 / 0 / 0
ORAL2013	<i>proboba / probůh</i>	53 / 0	16,13 / 0
ORAL2013	<i>pro boba</i>	7*	5,04
ORAL2013	<i>prokřistapána / propánajána / propánaboba</i>	0 / 1 / 0	0 / 0,3 / 0
ORAL2013	<i>pro křista pána / pro pána jána / pro pána boba</i>	0 / 0 / 0	0 / 0 / 0
ORTOFONv2	<i>proboba / probůh</i>	23 / 2	8,98 / 0,78
ORTOFONv2	<i>pro boba</i>	1*	0,39
ORTOFONv2	<i>pro křista pána / pro pána jána / pro pána boba</i>	2 / 0 / 0	0,78 / 0 / 0
ORTOFONv2	<i>prokřistapána / propánajána / propánaboba</i>	1 / 0 / 0	0,39 / 0 / 0

* (také) v adekvátní funkci předložkové substantivní vazby

Zdroj: vlastní

.. každopádně si myslím že by ženy se svými vědomostmi . @ . at už životními nebo že škol . neměly zůstat takzvaně na oet ale **probůh** už ne žádnou další Anežku Hodinovou Vzpurnou výstřelky . @ . byly všedny . @ . no (BMK)

mě zamkly v pokoji .. co si vůbec myslely .. že jsem .. jo a Pavel říkal no když #s jim utíkala .. tak . pr* . **prokřistapána** copak měli dělat .. sestřička tě bonila po baráku a nemohla tě najít .. tak ji zamkli .. (ORTOFONv2)

. a jo , vlastně , já sem jí slíbila tu kaši , **prokřistapána** . čekáš , až ti to vystydne? ne , už du . juž odcházím . bože , bože . no , podívej . eee . (ORAL2008)

lusev + pujova a samý placení že jo

jo právě
 lusev + pujova *to stojí peněz potom **ježkovy voči** tak*
to je jistý to je pe peněz vid'*
 lusev + pujova *ž* já jim do tobo mluvit nebudu **propánajána** --- nákyj*
 jo to jako jistý (ORAL2013)

voduna *no to víš no tak --- nevíš jaká budeš ty* (smích)
 rehega *doufám že takováde ne*
 voduna *(se smíchem) no to nevíš tak jako*
 rehega ***pro boha svatýho***
 voduna *(smích) no a tak jako* (ORAL2013)

Světlna R. *.. tak tam byla nějaká mladá holka myslím si že nějaká středoškolačka že se zaučovala ne? .. učila se*
 Růžena J. (pousmání) *jo .. ne tak to já bych asi nechtěla aby ta se zaučovala na mně **probůh***
 Světlna R. (smích) *no .. a právě .. přede mnou teda byla nějaká holka .. nejdřív brali z prstu to už ti brali krev z prstu .. taký? .. nebo ne?* (ORTOFONv2)

V zájmu automatizovatelné lemmatizace a snazšího vyhledávání vychází dvouúrovňový korpus ORTOFONv2 z ortografického přepisu; v zásadě proto není divu, že emocionální výrazy zaznamenává především analyticky. Jeho (popředložkové) komponenty jsou však tagovány jako substantiva. Tím se zcela stírá rozdíl mezi interjekcí a substantivním užitím, doloženým v témže korpusu kupř. sekvencí:

*...to jsou .. to je **pohled** prostě **pro boha** jako tak jo jakože jo takové jako přirovnání jakože ..*
 (ORTOFONv2)

Tuto pasáž, věnovanou změnám forem a funkcí prostředků mluvené komunikace, jsem uvedl příkladem spojení *ty vole / tyoe / tyve / tye*, užívaného mj. jako kontaktní partikule. Při sběru materiálu pro korpusy, na nichž jsme spolupracovali, jsem považoval za potřebné všechny podobné názny funkčních změn zaznamenávat. Poslední zmíněný korpus ORTOFONv2 eviduje 2630 dokladů formy *vole* a jedinou formu *voe*. Všechny ortograficky ztvárněné výrazy jsou tagovány jako substantiva. Tímto postupem je zcela zastřena funkční diference partikule, interjekce a jména (jako takového významně frekvenčně omezeného):

Zbyněk S. *hochu ale co mě překvapil třeba Jenda NP .. tam byl a .. ty jo ale já tomu Babišovi fandím ..*
 Norbert R. (smích) *kokote jeden **vole** že to tak řeknu ..*

6 Závěr

Pokusil jsem se shrnout několik zásadních, dosti různorodých momentů, jež jsme seznali za hodné pozornosti při tvorbě mluvených korpusů s pomocí studentů filologických oborů. Ti, nakonec stejně jako my sami jsme se při obstarávání zvukových záznamů neformálních komunikčních situací a jejich přepisu museli zabývat značně obecnými otázkami pragmlingvistickými a sociolingvistickými (povaha komunikační situace), obecně lingvistickými (sémioticky heterogenní povaha komunikátu, počítaje v to prolínání prostředků lingválních, paralingválních a nonlingválních), lexikologickými a morfosyntaktickými (delimitace jednotek v proudu řeči a zachycení jejich relevantní formy).

Porovnáním více korpusů mluveného jazyka se snažím ukázat dynamiku sledování určitých jevů a zisky, popř. ztráty plynoucí ze změny intence. Dílčí omluvou budiž, že omezení prostoru pro individuální snahu transkriptora zaznamenat co nejdělejší obraz zvukové formy je veden objektivizačním úsilím. Právě to bylo i pohnutkou k ústupu od interpunkce psaných textů a přechodu k interpunkci pauzové. Dalším zásadním motivem bude (resp. dlouhodobě je) strojové učení a automatické zpracování mluvené formy řeči.

Literatura

- Jana BÍLKOVÁ, 2021: Nepravá hypotaxe v spontánních mluvených projevech. *Štýl – komunikácia – kultúra*. Ur. Zuzana Popovičová-Sedláčková. Bratislava: Univerzita Komenského. 417–427.
- Jana BÍLKOVÁ, Jiří ZEMAN, 2021: Diskurzní marker no a v mluvené češtině. *Lingvistika – korpus – empirie*. Praha: Ústav pro jazyk český AV ČR, v. v. i., 191–198.
- Renata BLATNÁ, 2006: *Víceslovné předložky v současné češtině*. Praha: NLN, Nakladatelství Lidové noviny: Ústav Českého národního korpusu, 9–17.
- Miroslav GREPL, 1995: Příruční mluvnice češtiny. Praha: NLN, Nakladatelství Lidové noviny, 1995.
- Jana HOFFMANNOVÁ, Marie MIKULOVÁ, 2011: Korpusy mluvené češtiny a možnosti jejich využití pro poznání rozdílných "světů" mluvenosti a psanosti. *Korpusová lingvistika Praha 2011 - 2 Vězeň a výstavba korpusů*. Ur. František Čermák. Praha: NLN, Nakladatelství Lidové noviny. 78–92.
- Eliška KLEŇHOVÁ, 2010: Interjekce v českém jazykovém systému. Kvalifikační práce magisterská, Filozofická fakulta Univerzity Karlovy.
- Eliška KLEŇHOVÁ, 2012: Postavení a užívání interjekcí v současné češtině. *Naše řeč*, roč. 95, č. 5, 238–254.
- Zuzana KOMRSKOVÁ, Petra POUKAROVÁ, 2018: Kdo, kdy a proč skáče komu do řeči aneb překryvy ve spontánní mluvené češtině. *Korpus – gramatika – axiologie*, č. 17, 41–56.
- Zuzana KOMRSKOVÁ, Petra POUKAROVÁ, Martin HAVLÍK, 2019: Překryvy replik. *Syntax mluvené češtiny*. Ur. Jana Hoffmannová, Jiří Homoláč, Kamila Mrázková. Praha: Academia. 102–116.
- Zuzana LAUBEOVÁ, 2020: *Mluvenost v dialogické elektronické komunikaci*. Kvalifikační práce dizertační,

Filozofická fakulta Univerzity Karlovy.

Jan VOLÍN, 2015: Vztah mluvené a psané formy jazyka. *Mluvnice současné češtiny. 1, Jak se píše a jak se mluví*. Ed. Václav Cvrček. Praha: Univerzita Karlova, nakladatelství Karolinum. 44–47.

Miloslav VONDRÁČEK, 2008: Diskrétnost řečových jednotek. *Čeština v mluveném korpusu*. Ur. Marie Koprřivová, Martina Waclawičová. Praha: Nakladatelství Lidové noviny / Ústav českého národního korpusu. 255–261.

Miloslav VONDRÁČEK, 2013: Vlastnosti slova a slovní druhy. *Akademická gramatika spisovné češtiny*. Ur. František Šticha. Praha: Academia, 2013. 79–89.

Miloslav VONDRÁČEK, 2018: Druhy slov. *Velká akademická gramatika spisovné češtiny*. Ur. František Šticha. Praha: Academia, 2018. 35–110.

Slovník spisovného jazyka českého, 1989. 2. vyd. Praha: Academia.

Slovník spisovné češtiny pro školu a veřejnost, 1994. 2. vyd., opr. a dopl. Praha: Academia.

Internetové zdroje:

ASSČ: Akademický slovník současné češtiny. Dostupný z: <https://slovníkcestiny.cz/uvod.php>.
Poslední přístup 15. 6. 2023.

BMK: Brněnský mluvený korpus. Dostupný z: <https://wiki.korpus.cz/doku.php/cnk:bmk>. Poslední přístup 15. 6. 2023.

ORAL2006: Dostupný z: <https://wiki.korpus.cz/doku.php/cnk:oral2006>. Poslední přístup 15. 6. 2023.

ORAL2008: Dostupný z: <https://wiki.korpus.cz/doku.php/cnk:oral2008>. Poslední přístup 15. 6. 2023.

ORAL2013: Dostupný z: <https://wiki.korpus.cz/doku.php/cnk:oral2013>. Poslední přístup 15. 6. 2023.

ORTOFONv2: Dostupný z: <https://wiki.korpus.cz/doku.php/cnk:ortofon>. Poslední přístup 15. 6. 2023.

PMK: Pražský mluvený korpus. Dostupný z: <https://wiki.korpus.cz/doku.php/cnk:pmk>. Poslední přístup 15. 6. 2023.

PSJČ: *Přřruční slovník jazyka českého. Kartotéka lexikálního archivu (1911–1991)*. Dostupný z: <https://psjc.ujc.cas.cz/search.php?> Poslední přístup 15. 6. 2023.