

Korpuszépítés

Nyelv és informatika – Pécs, 2022/23 tavasz

5. óra

Simon Eszter – Vadász Noémi

2023. április 1.

1. NYTK-NerKor
2. KorKor
3. MoMa korpusz
4. Ómagyar Korpusz
5. **Uralic Languages Under the Influence Database**

NYTK-NerKor

NER: Named Entity Recognition

- a felügyelt gépi tanuláson alapuló rendszerek nagy mennyiségű gold standard adatot igényelnek → így van ez a magyar NER-ben is
- a gold standard adat előállítása drága és erőforrásigényes → kevés van, azok pedig erősen domainspecifikusak, korlátozottak méretükben és hozzáférhetőségükben

Magyar NER korpuszok:

- Szeged NER korpusz: gold, kb. 200.000, hírek
- Criminal NE korpusz: gold, kb. 500.000, hírek
- hunNERwiki korpusz: silver, kb. 19M, wiki
- NYTK-NerKor korpusz

Az NYTK-NerKor korpusz

- 1 millió szavas
- gold standard
- kiegyensúlyozott válogatás 5 domainből (kb. 200.000 szó/domain):
 - szépirodalom
 - jogi szövegek
 - vegyes webes szövegek
 - hírek
 - Wikipédia-cikkek
- gold standard morfológiai elemzés a korpusz egyötödén
- CC-BY-SA 4.0 licenc alatt elérhető:
<https://github.com/nytud/NYTK-NerKor>

- adatformátum: CoNLL-U Plus
- NE címkekészlet: CoNLL2002 (*PER*, *ORG*, *LOC*, *MISC*)
- NE címkeprefixálás: IOB2 (CoNLL2002)
- morfológiai elemzés: Universal Dependencies v2

A korpusz további tulajdonságai

- a korpusz előfeldolgozása az *emtsv*-vel történt
- további morfológiai elemzés az *emMorph* címkekészletében
- hivatalos train–devel–test vágás: 80%–10%–10% → kiegyensúlyozott: minden műfaj, forrás és morfológiai annotáltság ugyanilyen arányban van
- összehasonlításképpen: az OntoNotes 5.0 angol nyelvű része kb. 1,5 millió tokent tartalmaz

név: egyedi referenciával bír, a világ egy egyedi entitására utal

- *Kosztolányi Dezső*
- *Szilas Menti Mezőgazdasági Termelőszövetkezet*
- *United Nations Educational, Scientific and Cultural Organization*
- *Déli-Shetland-szk.*
- *IBM*
- *Kiss János altábornagy utca*
- *Műgyetem*
- *The Coca-Cola Co.*
- *Kovács Pistike*

Az annotálás alapelvei I.

- Csak tulajdonneveket annotálunk. Nem annotálunk olyan frázisokat, amelyek ugyan a világnak valamely egyedi részére utalnak, de nem tulajdonnévvvel. Például a *József Attila Gimnázium* annotálandó, de a szövegben szereplő az *a sul*i frázis nem, hiába derül ki a szövegből, hogy melyik iskolára utal.
- A nevek nem kompozicionálisak. Mivel a nevek jelölete nem a név részeinek a jelöletéből áll össze, ezért a neveket nem bonthatjuk részekre az annotálásakor. Például a *Kossuth Lajos utca* egy földrajzi névként jelölendő, hiába van benne egy személynév. Ebből az következik, hogy mindig a leghosszabb nevet (a legkülsőbbet) jelöljük a jelölhetők közül.
- Nem annotálunk egymást átfedő vagy egymásba ágyazott neveket. Vagyis minden annotációnak be kell fejeződnie, mielőtt egy másik elkezdődik.

Az annotálás alapelvei II.

- A *tag-for-meaning* elvét követjük. Vagyis egy nevet mindig az aktuális kontextusnak megfelelő referenciája alapján annotálunk.
- Ha az azonosított név ragozott formában szerepel a szövegben, a raggal együtt, a teljes alakot annotáljuk.
- A nevek képzett alakjait nem jelöljük. Nem annotálandók tehát az olyanok, mint *magyarországi, fideszes, petőfieskedő*.
- Ha a név összetétel előtagja, és az összetétel alaptagja köznév, például *Horn-kormány, Tilos Rádió-hallgatók, TA-vezérigazgató*, akkor nem annotálandók névként.
- A névhez nem tartozik hozzá az esetleg előtte álló névelő. Kivétel az az eset, amikor a határozott névelő része a névnek, például *The Hague, The Times*.

Az annotálás alapelvei III.

- A név rövidítése (akronim, mozaikszó, monogram) is névként annotálandó.
- A szöveghez a névannotálás során nem nyúlunk hozzá, vagyis nem javítjuk ki a helyesírási hibákat, nem vonunk egybe különírt szavakat, és nem választunk szét egybeírtakat. Ha valamilyen éktelen hibát látunk az aktuálisan címkézendő névvel kapcsolatban, akkor azt külön fel kell jegyezni. Ez a szépirodalmi szövegekre nem vonatkozik: ott az az elv, hogy amit a szerző leírt, az sérthetetlen.

- PERSON: Valós és kitalált személyek neve, becenevek, művésznevek, álnevek. Ide tartoznak a kisebb, kevésbé strukturált embercsoportok, közösségek is.
- ORGANIZATION: Olyan csoportok nevei, amelyek valamilyen szervezett struktúrával rendelkeznek, mint például intézmények, vállalatok, kormányzati hivatalok, sportcsapatok, múzeumok, egyetemek.
- LOCATION: Földrajzilag vagy politikailag definiált helyek nevei, úgymint városok, országok, hegyek, völgyek stb. Ide tartoznak az emberalkotta építmények is, mint a repterek, utak, gyárok, épületek stb.
- MISC: A felsorolt típusok egyikébe sem tartozó nevek.

2021-es megjelenése óta:

- sokan használják
- a Stanza (Qi et al., 2020) magyar nyelvű NER modulja
- a HuSpaCy (Orosz et al., 2022) magyar nyelvű NER modulja
- kiértékelések: kellően nagy és heterogén → a rajta tanított rendszerek rendelkeznek azzal az általánosító képességgel, amivel eddig nem
- további munkálatok

KorKor

- többretegű, kézzel annotált korpusz
- morfológiai címke és tő kétféle címkékészlettel, dependenciaelemzés, zéró igék és zérónévmások, anaforikus- és koreferenciakapcsolatok
- kétféle fájlformátum
- a korpusz, a munkafolyamat és a felhasznált eszközök elérhetőek itt:
https://github.com/vadno/korkor_pilot
- CC-BY-4.0. licenc

	dokumentum	token (<i>conllup</i>)	token (<i>xtsv</i>)
huwiki	62	16,739	18,262
globv	32	7,760	8,799
TOTAL	94	24,499	26,581

A munkafolyamat

1. szöveggyűjtés
2. *emtsv* elemzés (*emToken*, *emMorph*, and *emTag*)
3. formátumkonverzió
4. kézi ellenőrzés (Google Spreadsheets)
5. formátumkonverzió
6. *emtsv* elemzés (*emDep*)
7. formátumkonverzió (*emCoNLL*)
8. kézi ellenőrzés (WebAnno)
9. a zérólétigék és az elliptált igén kézi beillesztése (plain text editor)
10. zérónévmásbeszúrás (*emZero*)
11. névmási anafora beillesztése (saját szkript)
12. kézi ellenőrzés és koreferenciaannotálás (Google Spreadsheets)
13. formátumkonverzió

A zérónévmások beillesztése

- a magyar pro-drop nyelv, az ige alanya, tárgya és a birtok birtokosa lehet droppolva
- egy szabályalapú szkript illesztette be ezeket
- a szabályok a többi nyelvi annotáción (tő, morfológiai címke, dependenciaelemzés):
 - alany, ha a függőségi fában az igehez nem tartozott alany;
 - tárgy, ha egy határozott ragozású igehez nem tartozott tárgy a függőségi fában;
 - birtokos, ha egy birtokhoz nem tartozott birtokos a függőségi fában;
 - alany a ragozott vagy a ragozatlan infinitívushoz
- plusz ágak jelennek meg a függőségi fában
- a beillesztett névmások morfológiai jegyei (szám, személy) az ige vagy a tárgy morfológiai címkéje alapján kiszámolható
- 867 zéró alanyt, 101 zéró tárgyat és 379 zéró birtokost illesztettünk be

A névmási anaforikus kapcsolatok beillesztése

- egy szabályalapú szkript illesztette be
- a szkript megkeresi a személyes névmásokat
- a szabályok a szófaj, a morfológiai jegyek és a szintaktikai információk alapján működnek
- antecedenst keres a névmásoknak
- pl. ha az ige alanya zéró és az ige ragozása megegyezik az előző tagmondat igéével, akkor a zéró alany antecedense megegyezik az előző tagmondat alanyának antecedensével

Kézi javítás és koreferenciaannotáció

- 4 nyelvészhallgató ellenőrizte és javította a beillesztett zérónévmásokat és az anaforikus kapcsolatokat, valamint a annotálta a koreferenciakapcsolatokat
- Google Spreadsheets feltételes formázásokkal
- anafora típusok: személyes (**prs**), mutató (**dem**), kölcsönös (**recip**), visszaható (**refl**), vonatkozó (**rel**), birtokos (**poss**)

általános

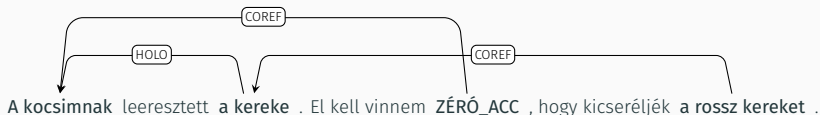
a Kínai Kommunista Párt egyik volt vezetője, akit hazaárulás miatt elítéltek

beszélő és címzett

A születésnap ajándékoknak is nagyon **örülünk**, ha **szeretnéd** támogatni a munkánkat, **küldj nekünk** adományt, vagy **vegyél** egyet az NSA-s karácsonyi **üdvözlőlapjaink** közül, amelyet a Creative Time-nál dolgozó **barátaink** terveztek.

Koreferencia és a szerteágazó kapcsolatok

- koreferenciatípusok: két elemnek ugyanaz a referense (**coref**), rész-egész kapcsolat áll fenn a két elem között (**holo**)
- szerteágazó kapcsolatok: a koreferenciakapcsolat csak testes elemek között lehet, a névmás antecedense lehet droppolt névmás is.



- A referens állapota megváltozik. Vajon a holttest koreferens az emberrel?

*Három hónap telt el az **újságíró házaspár**, Sagar Sarwar és felesége, Meherun Runi meggyilkolása óta. A **holttesteket** már exhumálták is, hogy megismételjék a boncolást.*

- split antecedens

***Papyrus** bátor és megmenti **Thèti-Chèri-t**. A két egymásra lelt **barát** küldetést kap az istenektől, hogy védelmezzék meg a fáraót.*

MoMa korpusz

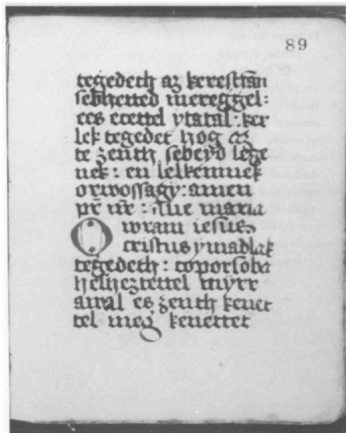
- moldvai magyar nyelvjárások
- Tánczos Vilmos kolozsvári néprajzkutató moldvai gyűjtései
- minden olyan moldvai települést felkeresett, ahol azt feltételezte, hogy tetten érhetők még addig nem dokumentált archaikus népi imádságok, így értékes anyag gyűlt fel a magyar ekevessé dokumentált változatairól
- előzmény: <https://nlp.nytud.hu/csango/>
- eddig nem létezett ezeknek a nyelvjárásoknak a kutatására alkalmas, megfelelő pontossággal lejegyzett és nyelvi annotációval ellátott korpusza
- jelenleg 63 ezer token (30 interjú 21 kutatópontról)
- augusztusra várhatóan 25 kutatópontról 60 interjú lesz

- lejegyzés a magyar egyezményes hangjelölés standardjával
- automatikusan egyszerűsített átirat
- a hangfelvételek össze vannak kapcsolva a szöveggel
- **normalizálás**
- morfológiai elemzés és tövesítés (az ómagyar elemző reszelgetésével)
- SketchEngine lekérdezőfelület
- nemcsak lekérdezőfelületen keresztül, hanem teljes elérés



Ómagyar Korpusz

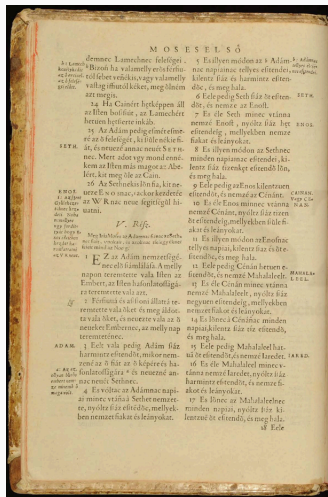
- 2009-2013, 2015-2018
- annotált korpusz, amely tartalmazza az összes fennmaradt ómagyar kori (896–1526) és néhány középmagyar kori (1526–1772) szövegemléket, és amely nyelvészeti releváns kérdésekre tud választ szolgáltatni
- 3,2 millió szövegszó: 47 ómagyar kódex, 24 rövidebb ómagyar szövegemlék, 244 misszilis és 5 középmagyar kori bibliafordítás
- <http://omagyarkorpusz.nytud.hu>



177
89r

tegedeth az keresztian
sebhettek mereggel :
ees ecettel ytatal : ker-
-lek tegedet hog az
5 te zenth sebeyd legé-
-nek : en lelkennek
orwossagy : amen
pr nr : Aue maria
O wram iesus
10 cristus ymadlak
tegedeth : coporsoba
helhezttel myrr-
-awal es zenth kenet-
-tel meg kénétet

Korai nyomtatványok, revideált változatok



MÓZES L. KÖNYVE 4. 5.

1

feleségeinek: Oh Hada és Cirila, hallgassatok szóra, Lámekh feleséget, halljátok beszédem; embert öltem, mert megsebzett; ifjat öltem, mert megütött.

24. Ha * hétszerez a bosszú Káinért, hetvenhétzerez az Lámekhéért.

Sally G. Fiske

25. Ádám pedig látni ismerte az ő feleségét, és az a nőle neki fiat, és nevezé annak nevét * Sétének: mert adott ágyamoz, ésenkem az Isten más magot Ábel helyett, ki megölte Káin. * *évek 5.3.*

26. Sétének is született fia, és nevezé annak nevét Enósnak. Akkor kezdődék segítségül hívni az Úrnak nevét.

Ádám nemzetsége a Sétő ágán;
a Noéig való pátriárkák
(N. 1. 1. Kéz. 1. 1. 1.)

5 Ez az Ádám nemzetségének könyve. A mely napon teremte Isten az embert, Isten képmé-

2. Férflővő * és asszonygyá teremé Éket. és mondá Éket és

nevezése az ő nevüket Ádámnak,
a mely napon teremtetének.

3. Elt vala pedig Ádám száz harmincz esztendő, és nemze fiait az ő képére és hasonlatosságára és nevezé annak nevét Sétnek.

4. És telének Ádám napjai, mi-
nekutánna *Séthet nemzette, nyolcz-
száz eszlendőre, és nemze fiaikat és

5. És Mós Ádám egész életének ideje kilenczszáz harmincz esztendő;

6. Éle pedig Séth szar öt esztendő, és nemzé Enóót.

7. És éle Seth, minckéntanna Enóst
nemzette, nyolczszáz hét esztendeig;
és nemze fiaikat és leányokat.

9. Éle pedig Én is kilencven esz-

10. És éle Enós, minekutánna Kénánt nemzette, nyolczszáz tízenöt

11. És lőn Énő egész életének

estoyé Kilescasas de Chalendo; en me-
hala.

12. Éle pedig Kénán hetven esztendő, és nemzé Mahabibélt.

13. És ele Kénán, melyekutánna Mahalálélt nemzette, nyolcvanáz negyven esztendőig; és nemze fiait és leányokat.

14. És lőn Kénán egész életének ideje kilenczszáz tíz esztendő; és meghalt.

15. Éle pedig Mahaláléi hatvanöt
százötöt, és nemzét Járódot.

na Járódet nemzette, nyolczszáz harmincz esztendeig, és nemze fiait és leányokat,

17. Es lőn Maharábel egész életének ideje nyolczszáz kilencvenöt vartendő; és meghala.

19. És die Járod, minekutánna
Énókhót nemzette, nyolczszáz esz-

20. És lőn János egész életének

21. Éle pedig Énökké hatvanöt esz-

22. Es járt * Énók az Istennel,
mínekutánna Methuséláht nemzette,
háromezer esztendőlt: és nemze Sám-

23. És lőn Énőkh egész életének
deje háromszáz hatvanöt esztendő.

24. És mivel Énőkh Istennel járt *
vala; eltűnök, mert Isten magához
vevő. * Zsolt. 115.

25. Éle pedig Metruselah száz
nyolczvanhét esztendő, és nemzé
Lámekhet.

28. Es de Mchuselah, mnek-
stánna Lánekhet nemzette, hét-
száz nyolcvankét esztendő; és
ezek a fiókai és leányai:

27. És látta Methuselah egész életének ideje kilencszáz hatvankilenc esztendő: és meghalt.

28. Éle pedig Lámekh száz nyolcz-
vankét esztendő, és nemze fiat.
29. És meverzé azt Noénak, mond-

rám: Ez vígasztal meg minket munkálkodásunkban s kerünk terhes fátadozásában e földön, melyet meg-

30. És éle Lámeké, minékutánna
Noét nemzette, ötszár kilenczevenöt
évet élte: és nemze fiait, és lef-

ryokai.

Érvek a szerkesztett kiadások...

...mellett

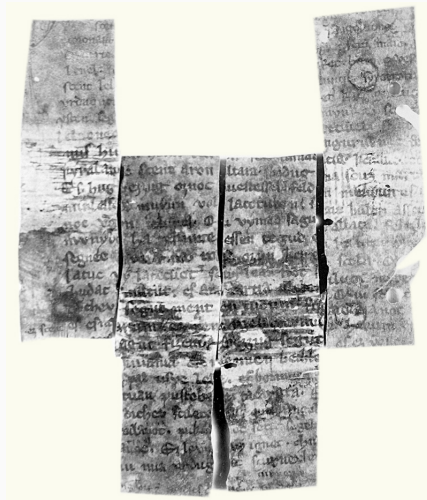
- könnyen elérhetőek
- a szerkesztői döntéseket már meghozták
- könnyebben szkennelhető vagy begépelhető

...ellen

- a szerkesztést nem nyelvészek végezték
- a meghozott nyelvészeti döntések nincsenek kellően dokumentálva
- nem alkalmasak további nyelvészeti vizsgálatokra
- copyright

Kiút a dilemmából: szerkesztett kiadásokból kiindulni, de mindent ellenőrizni az eredeti verzióban.

Rossz fizikai állapot



Speciális karakterek

52 latin alapkarakter

42 diakritikus jel

10 szám

15 speciális karakter

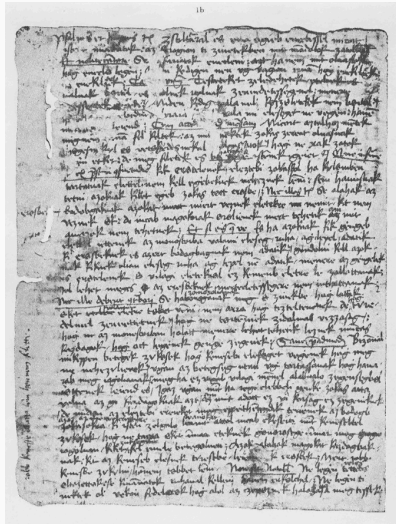
34 szövegtagoló és egyéb jel

3 görög betű

Összesen: 156 karakter + ezek kombinációi

í J Ñ R ā ē ġ ħ ñ ó ů v | ~ ¶ ß ě Ÿ ě [Θ] 3 ö œ 3 ^{ra} Γ α ħ r ŵ † 7 : p p ² ³ ⁴ ⁵ ⁶ ⁷ ⁸ ⁹

Betoldások, javítások, etc.



Heterogén helyesírás, normalizálás

a már nem létező jelenségeket megőrizni

ýsa	pur	es	chomuv	uogmuc
isa,	por	és	hamu	vagyunk
lata	q	napat	fèkette	
látá	ő	napát	fekette	

a helyesírási esetlegességeket eltörölni

meden ~ minden ~ mendun ~ mendē ~ mendē ~ miden ~
minden ~ mynden ~ mýnden ~ mýndē ~ mýden ~ mýnden ~
mýndew ~ mýnden ~ mýndon ~ mēden ~ mēdèn ~ mēdē ~
mēdèn ~ mēnden ~ mēden → minden

Megoldás: az eredeti és a normalizált verziót is tartalmazza a korpusz.

az elérhető elemzők a modern nyelvállapotra készültek

- a régi nyelvállapot elemzésére alkalmas eszközöket kell fejleszteni
- az adaptáció nem triviális
- több kézi ellenőrzést igényel

Uralic Languages Under the Influence Database

- *Az uráli nyelvek mondattanának változása aszimmetrikus kontaktushelyzetben*
- 2016. február – 2017. július
- MTA Nyelvtudományi Intézet
- interdiszciplináris csapat: kutatók a finnugor, a nyelvtechnológiai és az elméleti nyelvészeti osztályról
- projektvezető: É. Kiss Katalin
- a vizsgált nyelvek: udmurt, tundrai nyenyec, színjai és szurguti hanti
- célja egy pilot adatbázis építése volt:
 - 4000 token/kor/nyelv
 - IPA-átirat
 - teljes morfológiai elemzés
 - angol fordítás

UDMURT:

- írott (5) kategória: napi szinten használják, és létezik egy sztenderd irodalmi változata, de az nem annyira terjedt el
- Udmurtia egyik hivatalos nyelve
- vannak udmurt blogok, napi szinten keletkezik elektronikus udmurt szöveg
- udmurt Wikipédia (1000+ cikk)

TÖBBIEK:

- veszélyeztetett (6b) kategória: csak informális körben használják, alacsony presztízs
- nincs Wikipédia
- nincs napi sajtó, nem keletkezik elektronikus szöveg

- **elsődleges adatok:** olyan kommunikációs eseményekből származó nyelvi adatok, amelyek a hétköznapi nyelvhasználatot tükrözik
- **teljességre törekvés:**
 - minél több társadalmi osztályt, kort, nemet, műfajt és dialektust reprezentáljunk
 - metaadatok
 - az eredeti felvétel megőrzése, hogy a lejegyzések ellenőrizhetők legyenek
- **egységesség és összehasonlíthatóság:** nemzetközi sztenderdek
 - Unicode
 - IPA
 - Leipzig Glossing Rules
 - UTF-8 kódolású tsv fájlok

A korpusz felépítése

YRK Hajdú:	jā	mīdaxana	amkerta	jaŋkūwi
YRK Mus:	ja	midaxana	amkerta	jaŋkuwi
YRK IPA:	ja	mi:daxana	ămkerta	jăŋkuwi
YRK cirill:	я	мыдахана	амкэрта	яңкувы
lemma:	я	мы	ңамгэ	яңгось
szófaj:	N	Ptcp	Pron.neg	V
glossza:	earth	create.IPFV.PTCP.LOC	what.CONC	neg.EX.INFER

ENG: when the earth was created, there was nothing

GER: zur zeit der erschaffung der erde gab es nichts

HUN: a föld teremtésének idején nem volt semmi

<http://archive.nytud.hu/oszt/elmnyelv/urali/adatbazisok.html>

Simon (2017)

Simon and Mus (2017)

References

- Orosz, Gy., Szántó, Zs., Berkecz, P., Szabó, G., and Farkas, R. (2022). HuSpaCy: an industrial-strength Hungarian natural language processing toolkit. In *XVIII. Magyar Számítógépes Nyelvészeti Konferencia*.
- Qi, P., Zhang, Y., Zhang, Y., Bolton, J., and Manning, C. D. (2020). Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.
- Simon, E. (2017). Négy hatás alatt álló nyelv – korpuszépítés kis uráli nyelvekre. In Vincze, V., editor, *XIII. Magyar Számítógépes Nyelvészeti Konferencia (MSZNY2017)*, page 263–274, Szeged.

Simon, E. and Mus, N. (2017). Languages under the influence: Building a database of uralic languages. In M. Tyers, F., Riessler, M., Pirinen, T. A., and Trosterud, T., editors, *Proceedings of the Third International Workshop on Computational Linguistics for Uralic Languages*, page 10–24, Saint Petersburg.