

# THE EHR LANGUAGE GARDEN

Leveraging Variability in Health Documentation

Denis Newman-Griffis

*NIH Clinical Center / University of Pittsburgh*



September 1, 2020



National Institutes of Health  
*Clinical Center*





Radiology

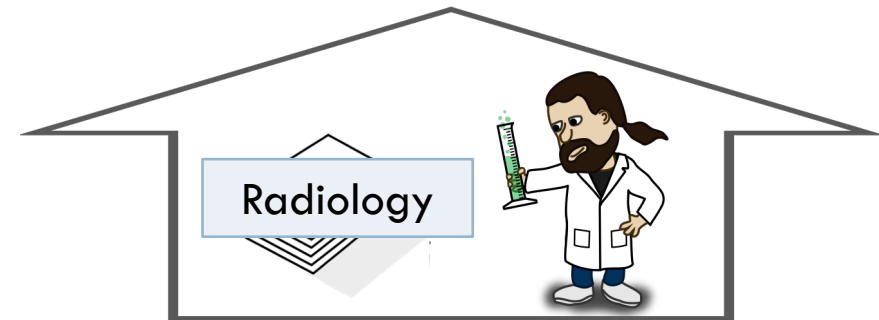
Pharmacy

Nursing

Discharge Summaries

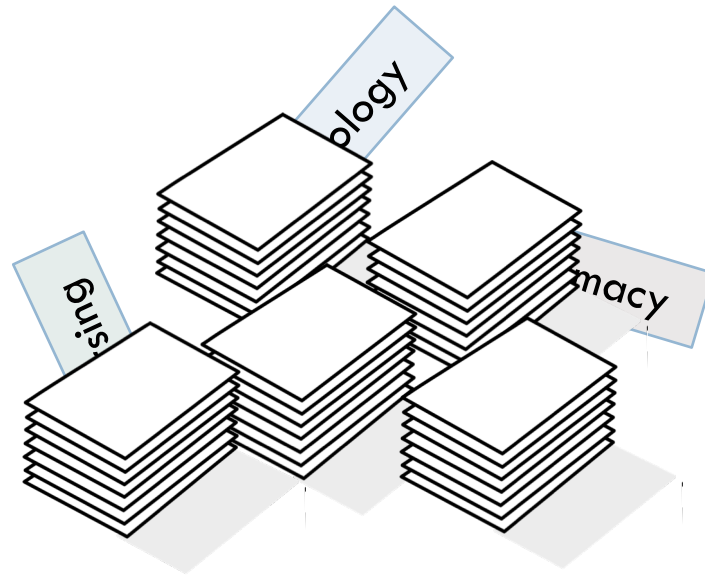
# Research → Practice

4



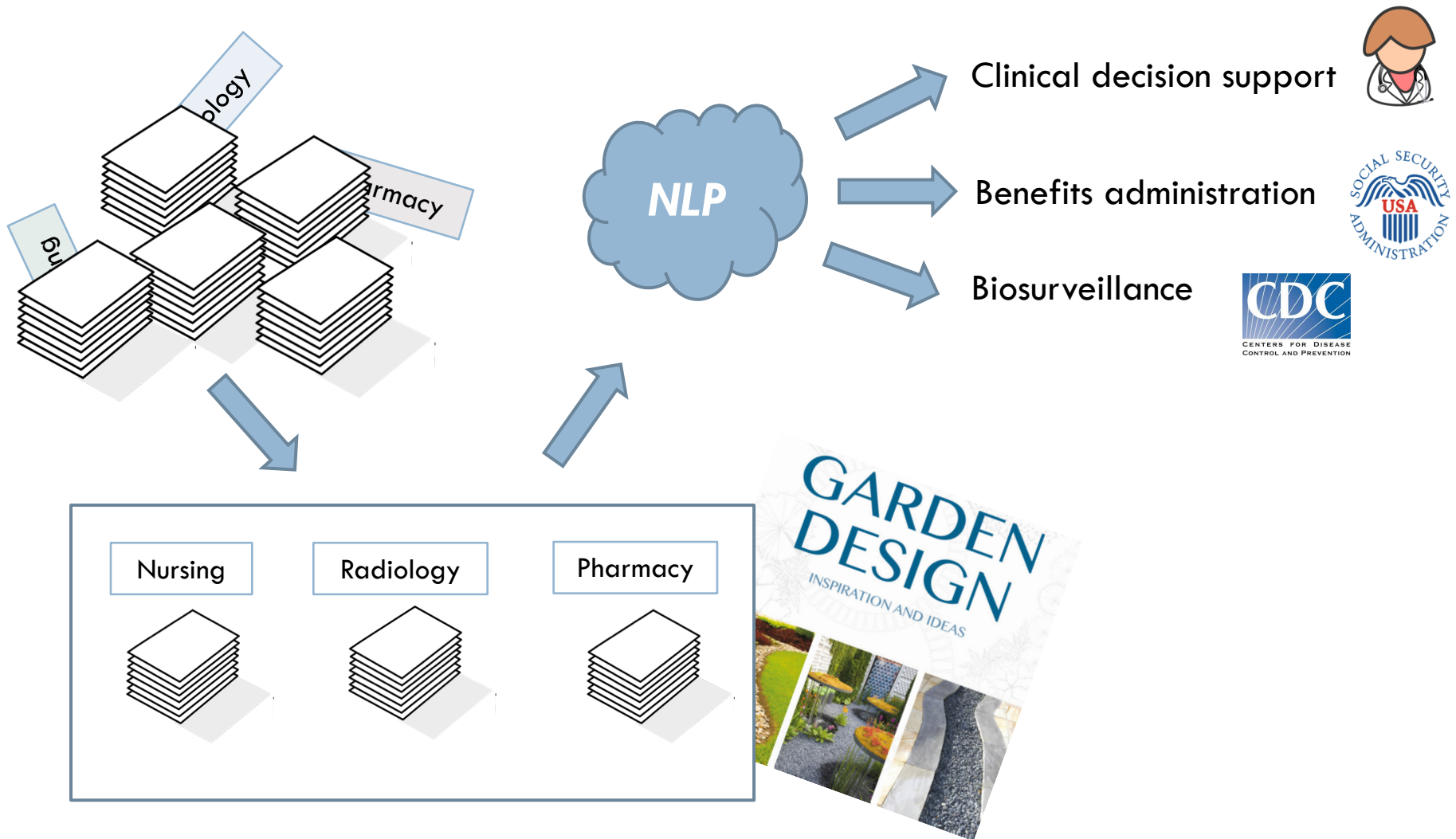
# Research → Practice

5



# Sublanguage: the secret sauce

6



# ERH data variability at scale

7

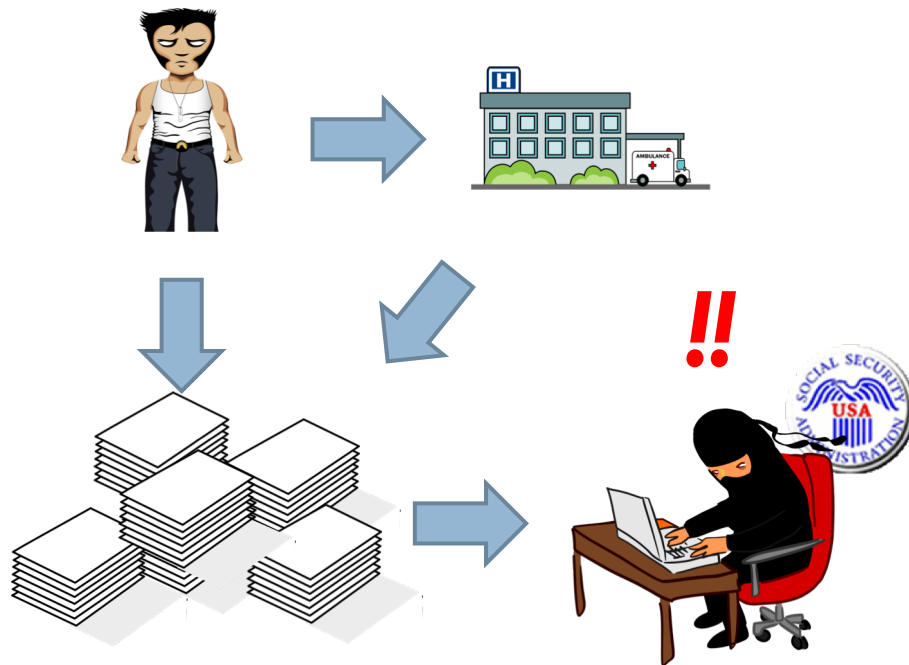
**Content** – what do the records say?

**Form** – how do they say it?

**Structure** – what are the pieces?

# Context: SSA disability programs

8



- ✓ National data
- ✓ All providers/EHRs
- ✓ Unreliable metadata



# “Defining” disability

9

## **Medical conditions**

- High mortality conditions
- Medical listings  
(business rules)

## **Functional limitations**

- Ability to perform work-related activities
- Substantial Gainful Employment

Need NLP that can handle both!

# Planting the garden: findings

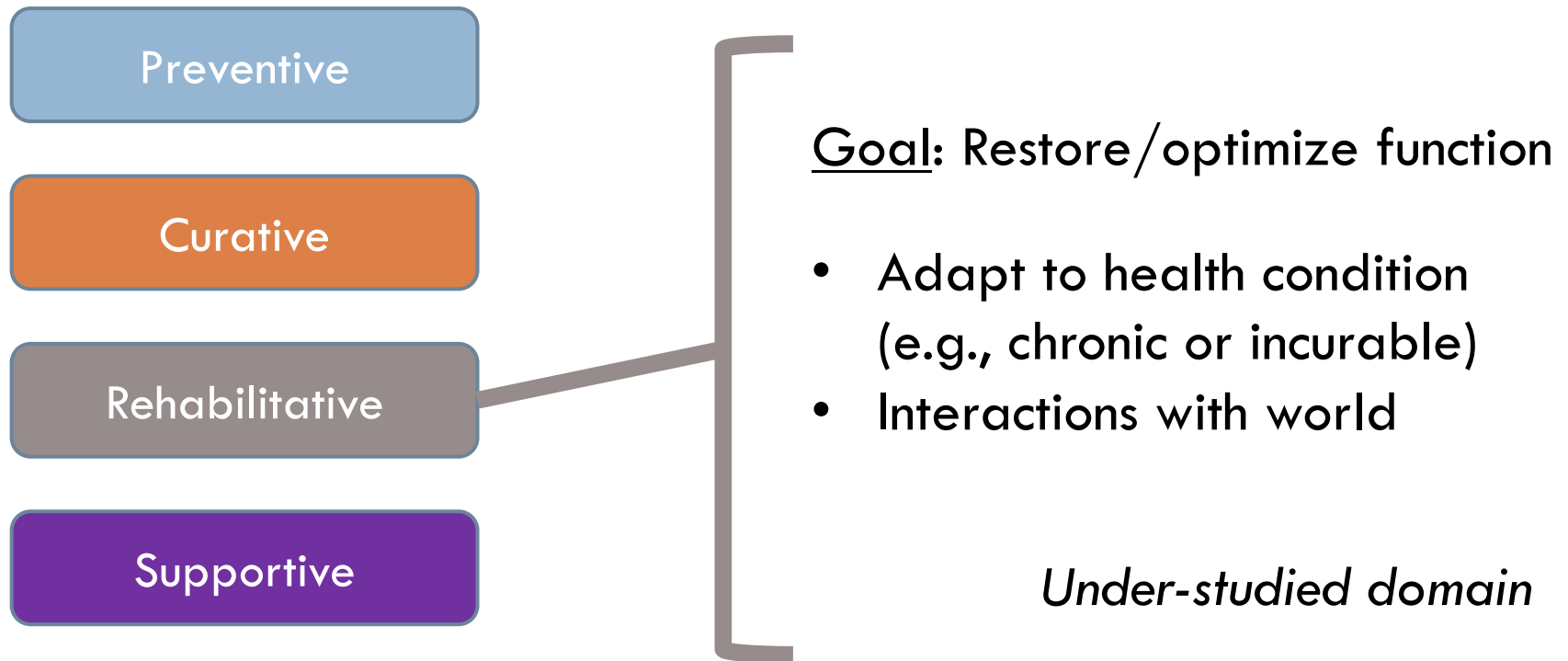
10

**Content** – Rehabilitation medicine as a sublanguage

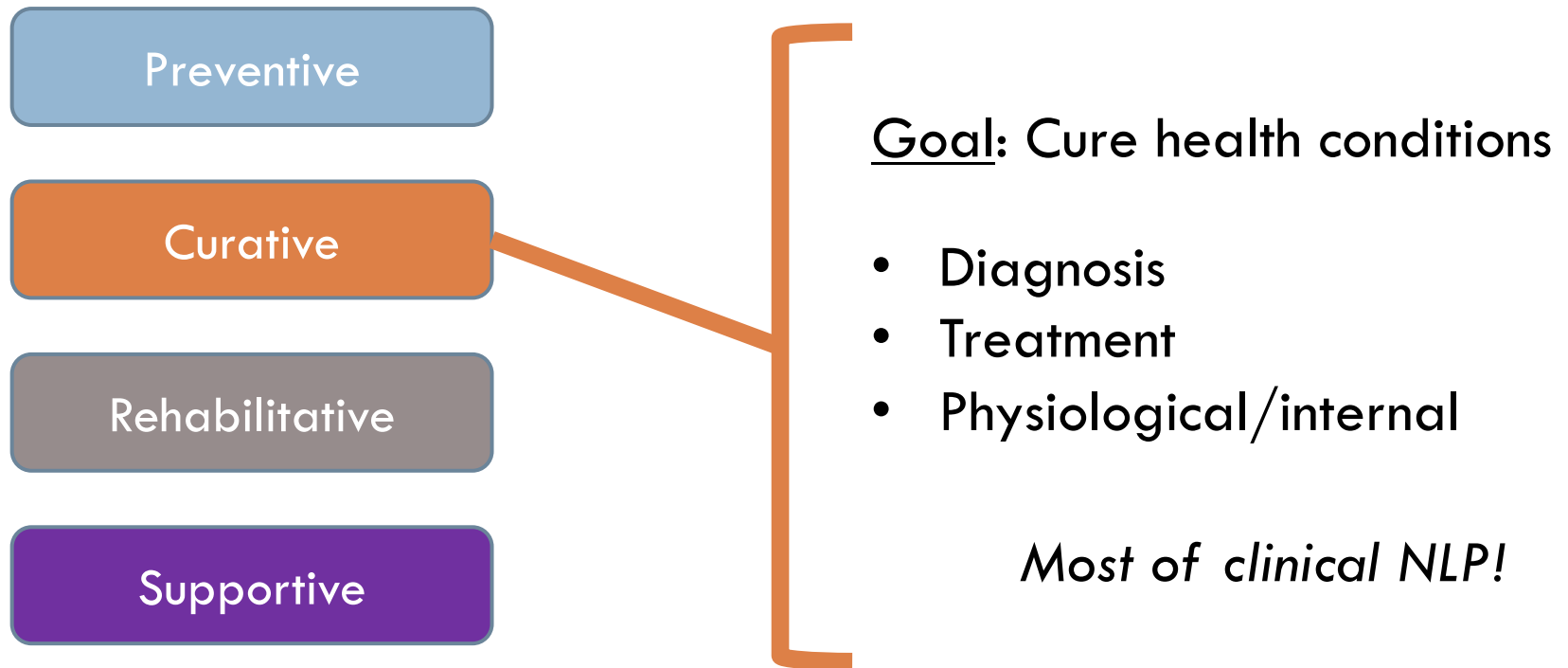
**Form**

**Structure**

# Health strategies – Rehabilitation



# Health strategies – Curative



# Multi-institution data

13



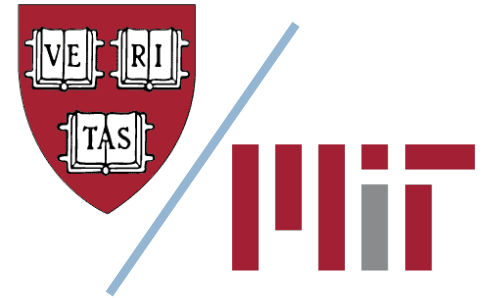
## BTRIS

- 155K records
- Research patients
- 130 doctypes



## OSUMC

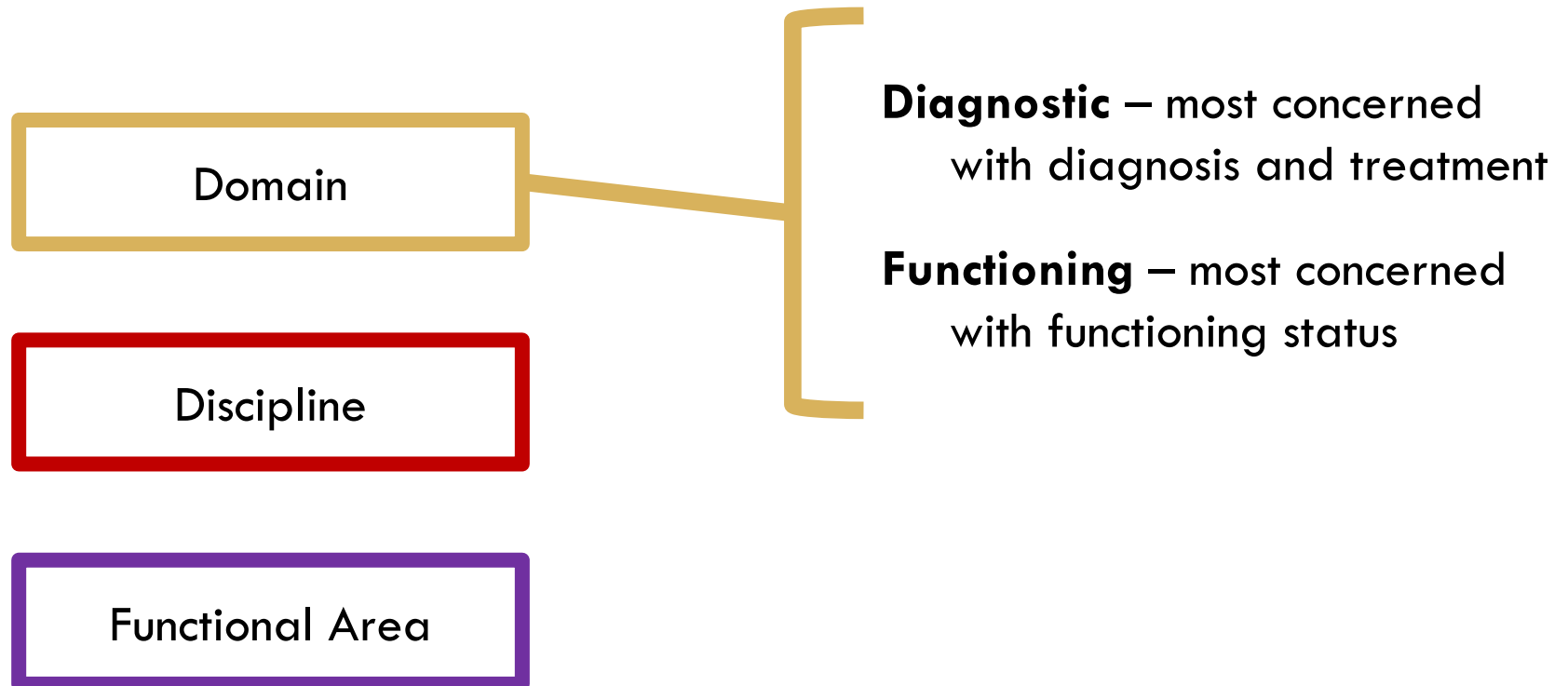
- 418K records
- Chronic diseases
- 43 doctypes



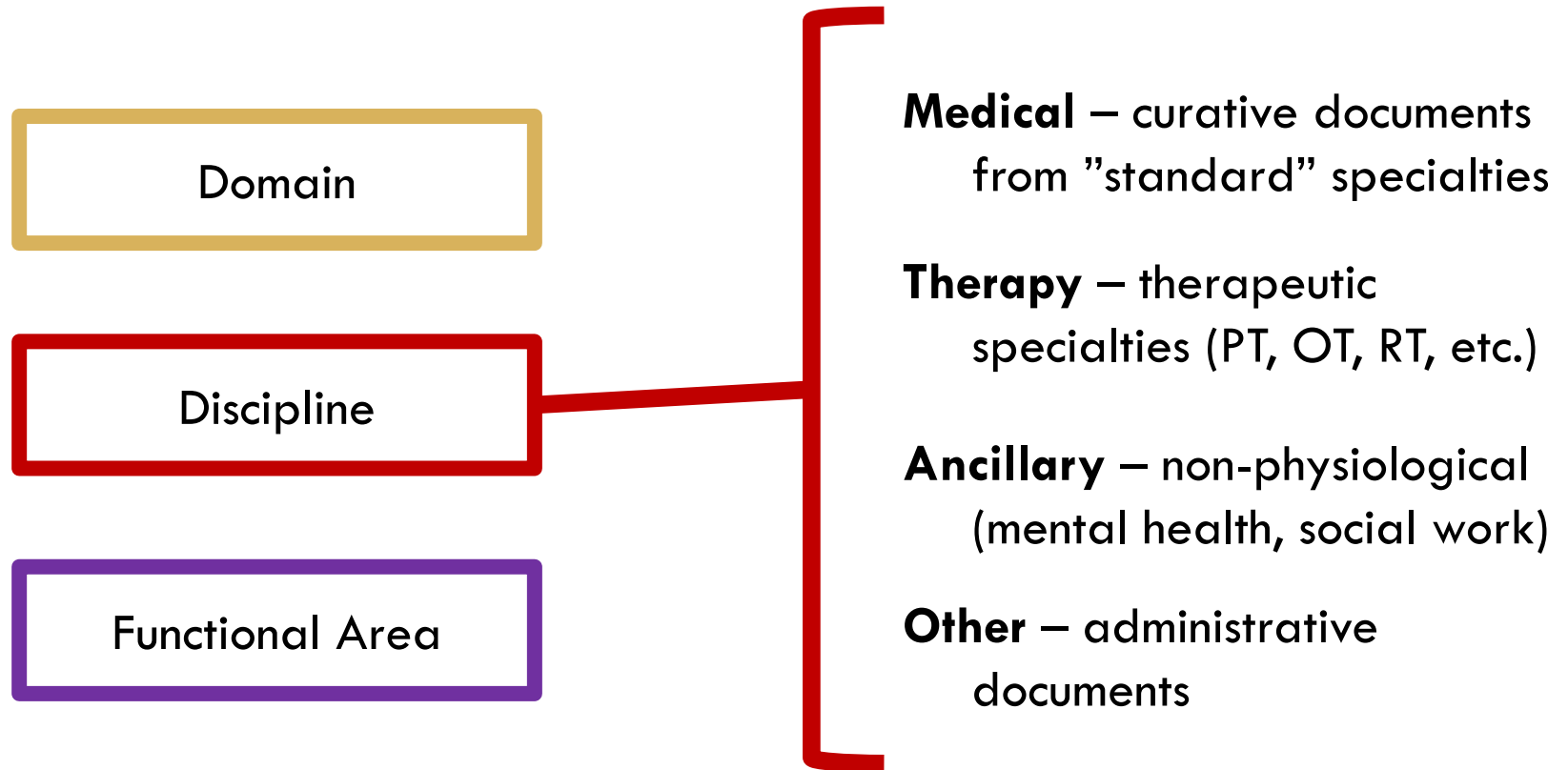
## MIMIC-III

- 2M records
- ICU admissions
- 25 doctypes\*

# Data classifications



# Data classifications



# Data classifications

Domain

Discipline

Functional Area

**PT** – physical therapy

**OT** – occupational therapy

**RT** – recreational therapy

**SLP** – speech/language  
pathology

**Psych** – psychological/iatric

**Neuro** – neurological

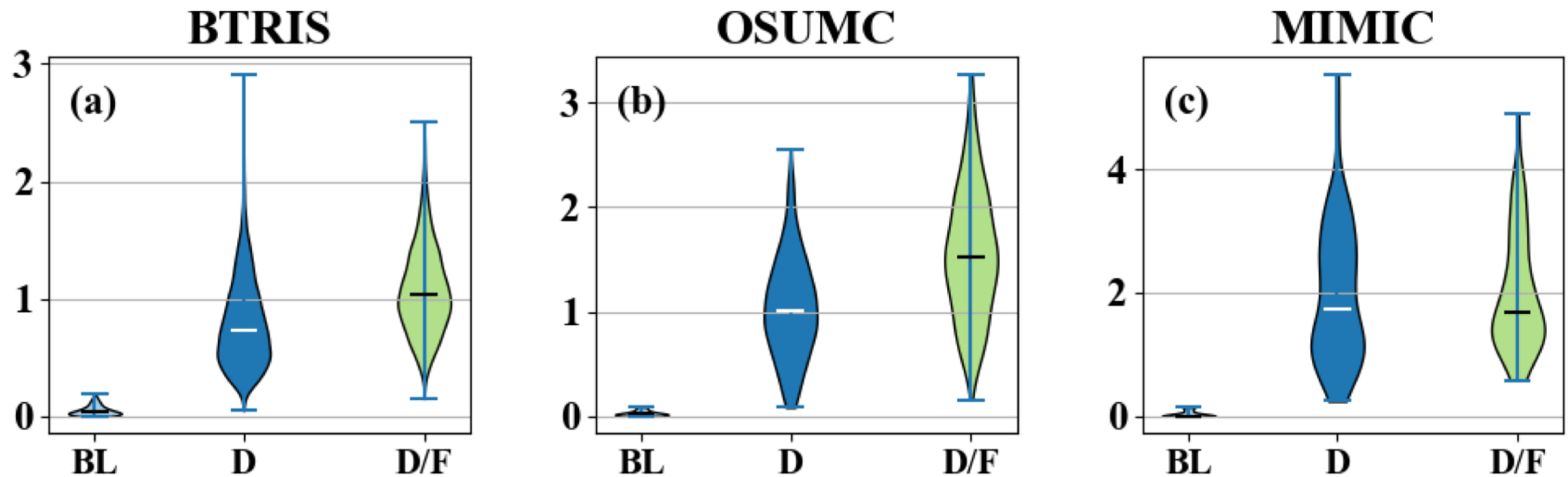
**SW** – social work

**General** – catchall bucket



# Rehab medicine vocabulary is distinct

17



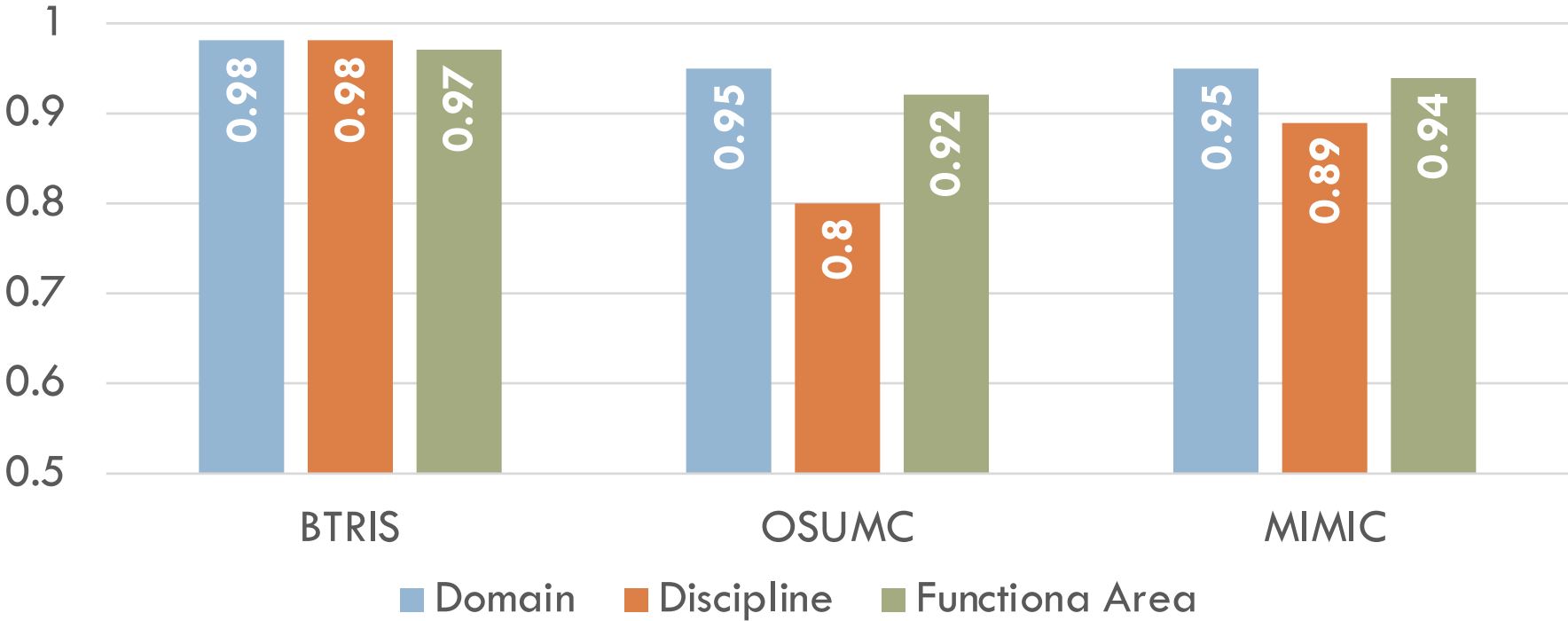
BL = Variance within document types

D = Variance between doctypes in Diagnostic

D/F = Variance between Diagnostic doctypes and Functioning doctypes

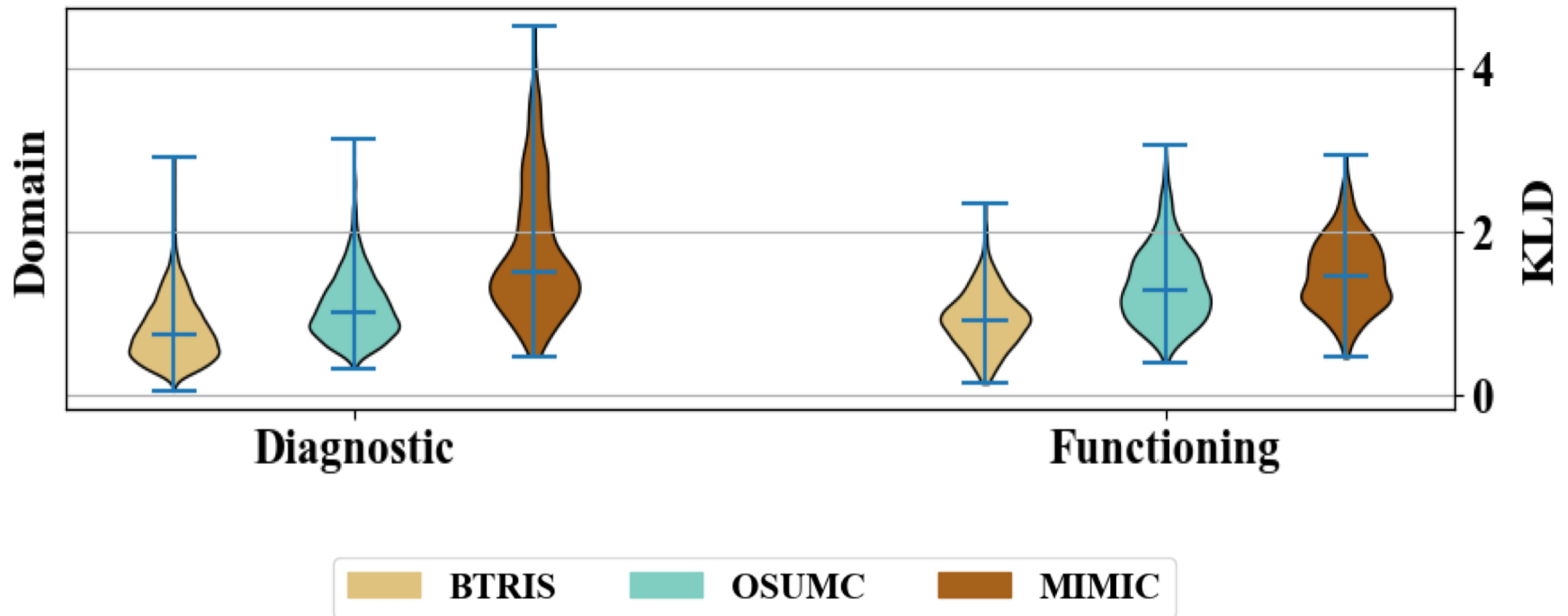
# Rehab medicine vocabulary is distinct

Document classification accuracy (unigram features)



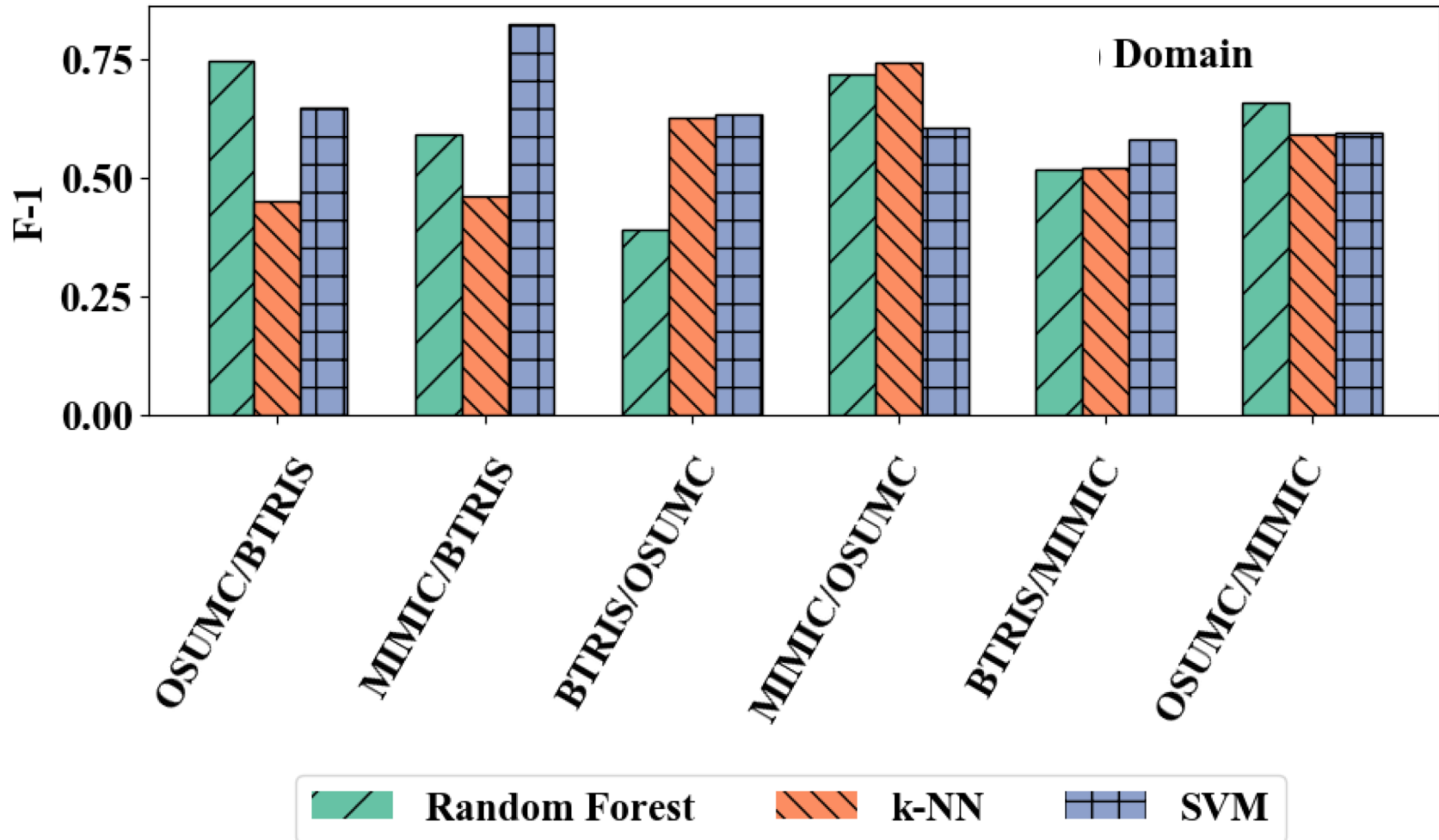
# Significant differences across institutions

19



# Significant differences across institutions

20



# Different structure of information

21

## i2b2

Ejection fraction: 90%  
Lab creatinine: 3 mg/dL

There has been removal  
of [a swan-ganz  
catheter]<sub>Treatment</sub> and  
placement of [a right  
internal jugular  
vascular  
catheter]<sub>Treatment</sub>.

## Rehab data

Pt 45 yr old tech worker,  
sedentary activity but  
hikes on weekends.

[Ambulation: 4]<sub>Mobility</sub>

Observations:

Pt is weight bearing: [she  
ambulates independently  
w/o use of assistive  
device]<sub>Mobility</sub>. Limited to  
very brief examination.

# Planting the garden: findings

22

## **Content**

**Form – Differences in clinical concept usage**

## **Structure**

**D N-G**, E Fosler-Lussier. “Writing habits and telltale neighbors: analyzing clinical concept usage patterns with sublanguage embeddings.” *LOUHI*, 2019.

# Characterizing document types

23

Document/section structural patterns inform meaning

- Field names vs observations
- Temporality (future/past/recurrent)
- Perceived importance (e.g. Chief Complaint)

Document types change priors for disambiguation

- “Depression” in Psychiatric Consult vs GE Exam

Discharge summaries  
!=  
Nursing notes

# Conceptual vs lexical analysis

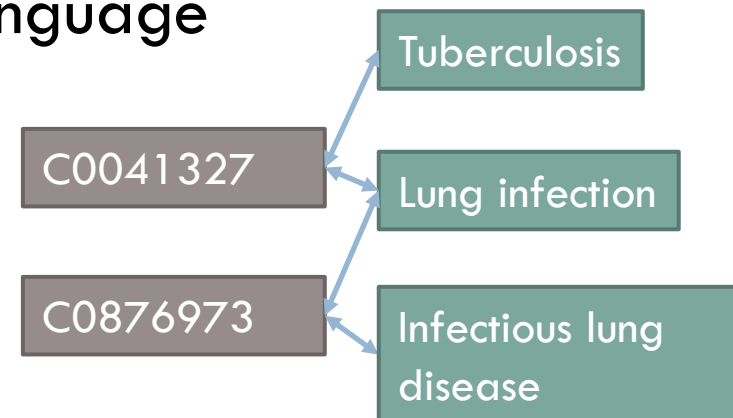
24

Prior work used lexical content to describe clinical sublanguages

- Feldman et al, 2016
- Grön et al, 2019

*Concepts* (symptoms, diseases, procedures, etc) are stock in trade of clinical language

- Multiple surface forms
- Ambiguity (“Cold”)





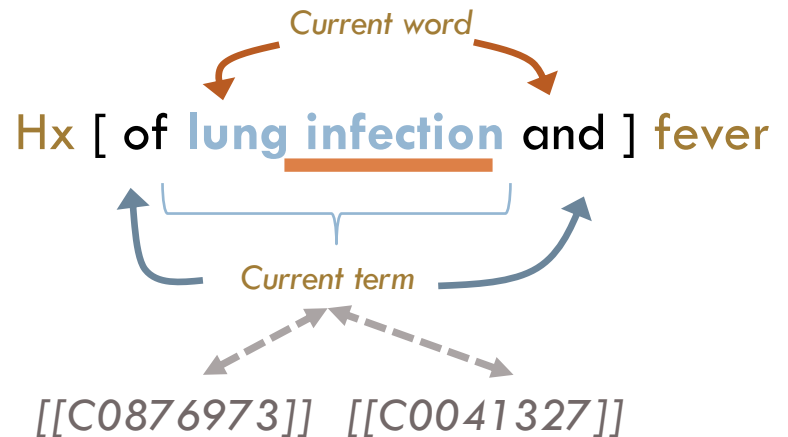
# Learning concept embeddings: JET

25

- Train word/term/concept embeddings jointly
- Distant supervision using known terminology
- Noisy, but good quality

Terminologies: SNOMED CT, LOINC

Data: MIMIC-III 

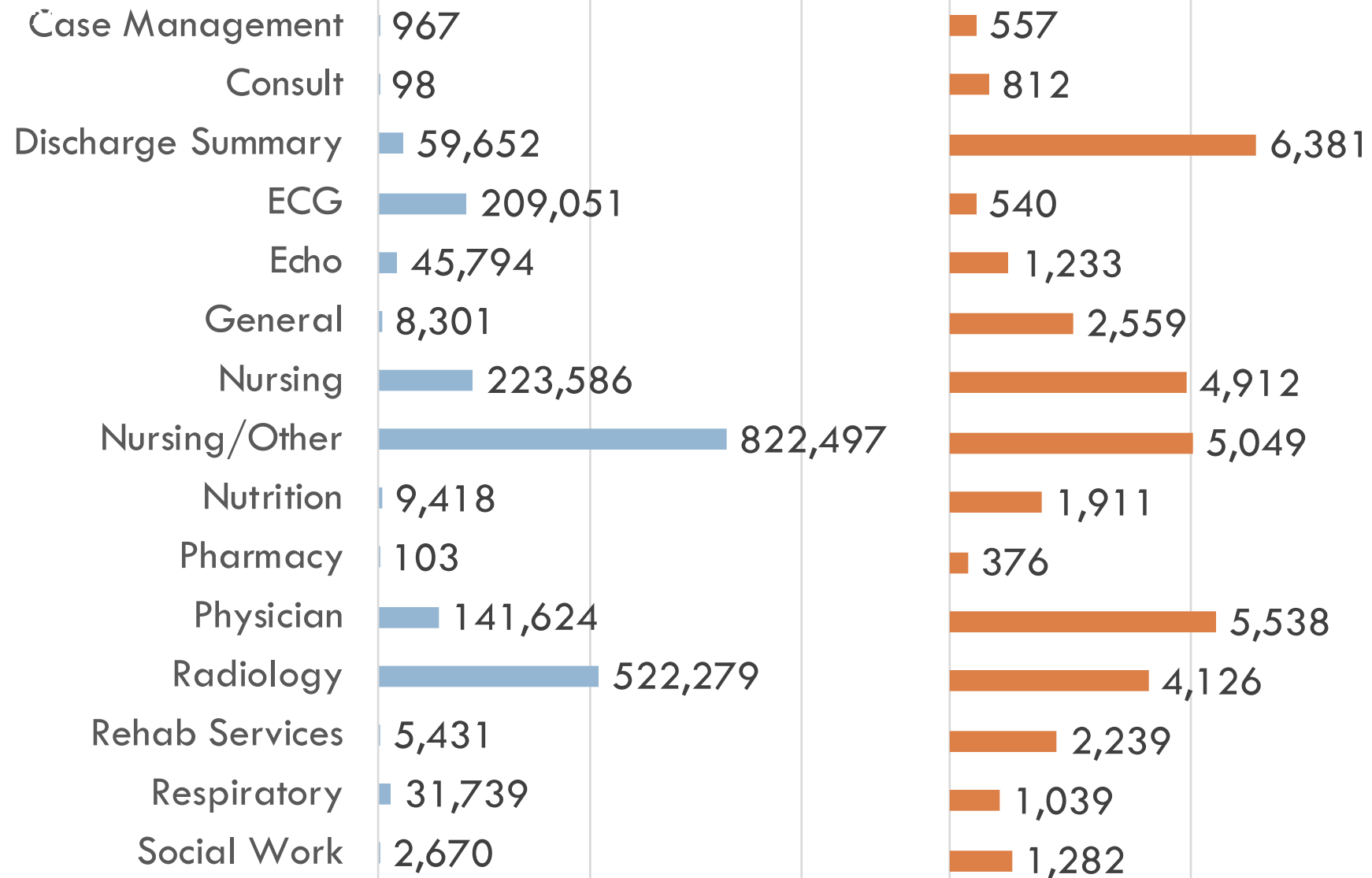


## # Documents

0 500,000 1,000,000

## # Concepts

0 5000 10000



# Measuring concept usage similarity

27

- ❑ Measured by overlap of nearest neighbor sets
- ❑ Similarity metric in  $[0,1]$
- ❑ Compare inter-type overlaps to intra-type overlaps

*Neighbors of Onion*

Set A	Set B
Cucumber	Squash
Squash	Pumpkin
Beans	Pasta
Green	Beans
Pasta	Cheese

# Inter-type similarity is significantly lower than intra-type

28

Case Management	0.75	0.01	0.01	0.00	0.00	0.00	0.01
Discharge Summary	0.01	0.67	0.24	0.32	0.00	0.34	0.33
Echo	0.01	0.24	0.65	0.13	0.00	0.36	0.40
Nursing/Other	0.00	0.32	0.13	0.60	0.00	0.27	0.31
Nutrition	0.00	0.00	0.00	0.00	0.73	0.01	0.00
Physician	0.00	0.34	0.36	0.27	0.01	0.57	0.26
Radiology	0.01	0.33	0.40	0.31	0.00	0.26	0.63
	Case Management	Discharge Summary	Echo	Nursing/Other	Nutrition	Physician	Radiology

# Nearest neighbors: Diabetes Mellitus (C0011849)

Discharge Summary	Nursing/Other	Radiology
Diabetes (C0011847)	Gestational Diabetes (C0085207)	Poorly controlled (C3853134)
Type 2 (C0441730)	A2 immunologic symbol (C1443036)	Insulin (C0021641)
Type 1 (C0441729)	Diabetes Mellitus, Insulin-Dependent (C0011854)	Diabetes Mellitus, Insulin-Dependent (C0011854)
Gestational Diabetes (C0085207)	Factor V (C0015498)	Diabetes Mellitus, Non-Insulin-Dependent (C0011860)
Diabetes Mellitus, Insulin-Dependent (C0011854)	A1 immunologic symbol (C1443035)	Stage level 5 (C0441777)

Strings: "diabetes mellitus",  
"diabetes mellitus dm"

# Nearest neighbors: Mental state (C0278060)

Discharge Summary	Echo	Radiology
Coherent (C4068804)	Donor [LOINC] (C3263710)	Mental status changes (C0856054)
Confusion (C0009676)	Donor person (C0013018)	Abnormal mental state (C0278061)
Respiratory status [LOINC] (C2598168)	Respiratory arrest (C0162297)	Level of consciousness (C0234425)
Respiratory status (C1998827)	Organ donor [LOINC] (C1716004)	Level of consciousness [LOINC] (C4050479)
Abnormal mental state (C0278061)	Swallowing (C4281783)	Mississippi (C0026221)

Strings: "mental status",  
"mental state"

# Embeddings pick up template patterns

## *Mental status* in Echo notes

PATIENT/TEST INFORMATION  
Indication: Pt presents with  
reduced mental status

PATIENT/TEST INFORMATION  
Indication: Pt presents in  
vegetative state, consider for  
organ donation

# Planting the garden: findings

32

**Content**

**Form**

**Structure** – Structural text features capture format and content variation



# Sources of variability in SSA data

33

## **Document Source**

- ❑ SSA Consultative Exams, CCDA documents from EHR, VA data, scanned notes

## **Content Types**

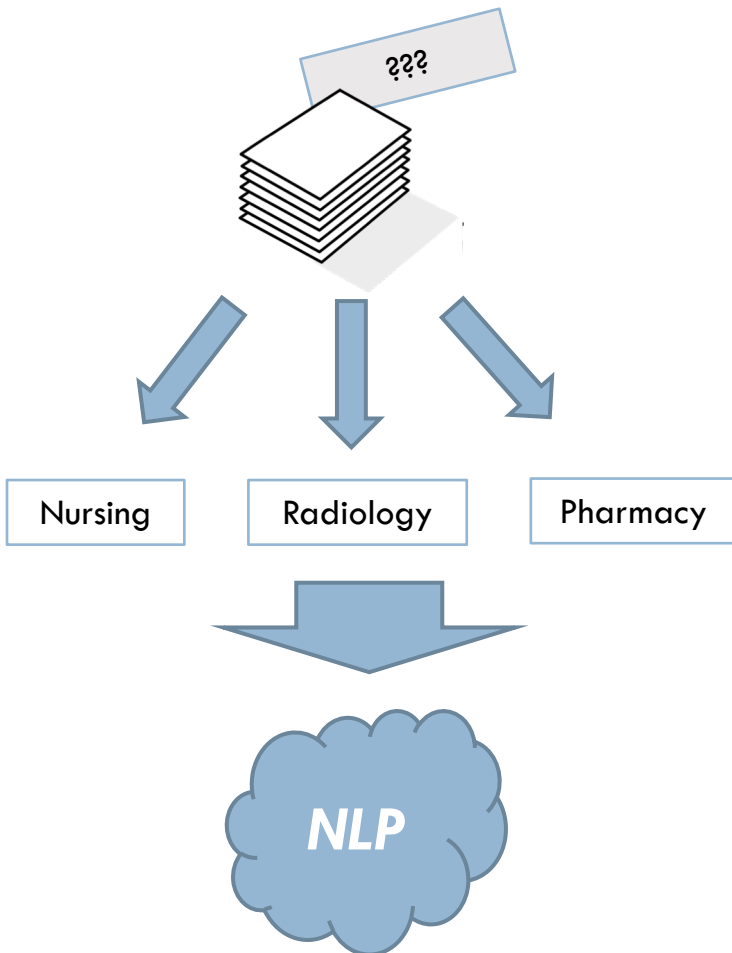
- ❑ SOAP notes, radiology reports, labs, surveys

## **Formatted Structure**

- ❑ Headers/footers, columns, section names, checkboxes

# Classify early, process better

34



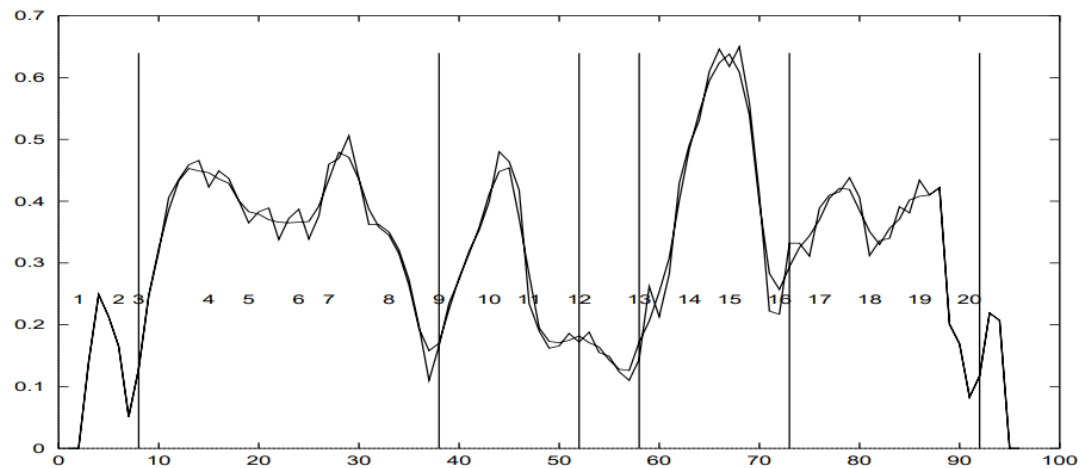
- 70K documents
- Disability claimants from 5 states
- Unreliable doctypes

# Page-level Features

- Number of Characters, Words, Lines, Sentences
- Number of Punctuation, Delimiters
- Number of Section Names, Section Zones, Nested Sections
- Number of Slot Values, Slot Names, Slot Value Values
- Number of Check Boxes (this wasn't actually working as it turns out)
- Number of Tables
- Number of Lists, List Elements
- Number of Questions
- "Text Tiling" Vector fingerprint (2 numbers)

# Related Work: Text Tiling

Marti Hearst (1994): Using word sequences to build a signal to indicate topic/paragraph shifts.



**Figure 6**

Results of the block similarity algorithm on the *Stargazer* text with  $k$  set to 10 and the loose boundary cutoff limit. Both the smoothed and unsmoothed plot are shown. Internal numbers indicate paragraph numbers, x-axis indicates token-sequence gap number, y-axis indicates similarity between blocks centered at the corresponding token-sequence gap. Vertical lines indicate boundaries chosen by the algorithm; for example, the leftmost vertical line represents a boundary after paragraph 3. Note how these align with the boundary gaps of Figure 5 above.

Marti A. Hearst, Multi-Paragraph Segmentation of Expository Text. *Proceedings of the 32nd Meeting of the Association for Computational Linguistics*, Los Cruces, NM, June, 1994.

# Page-level PCA

37

## Input

- ❑ Structural features for each page

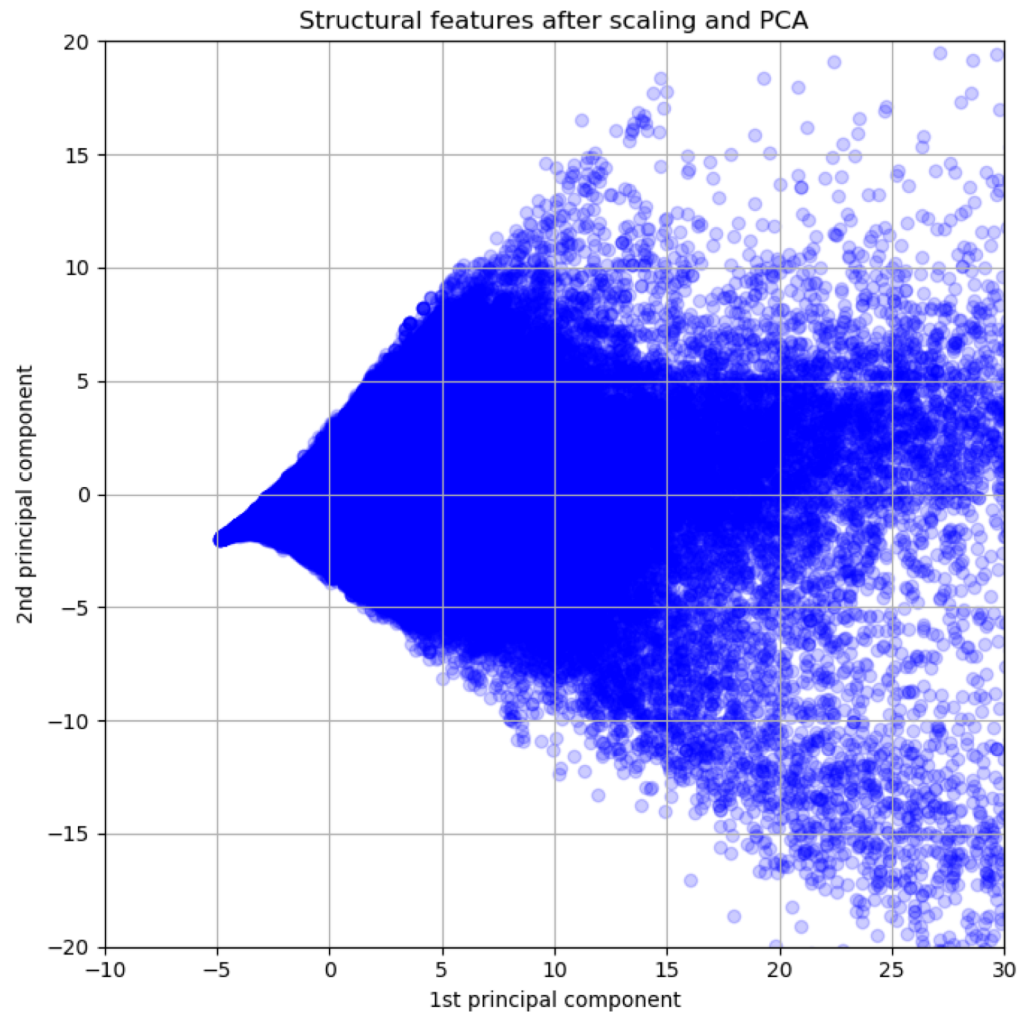
## Output

- ❑ Orthogonal transform into a set of principal components
- ❑ Dimensionality reduction and variance identification

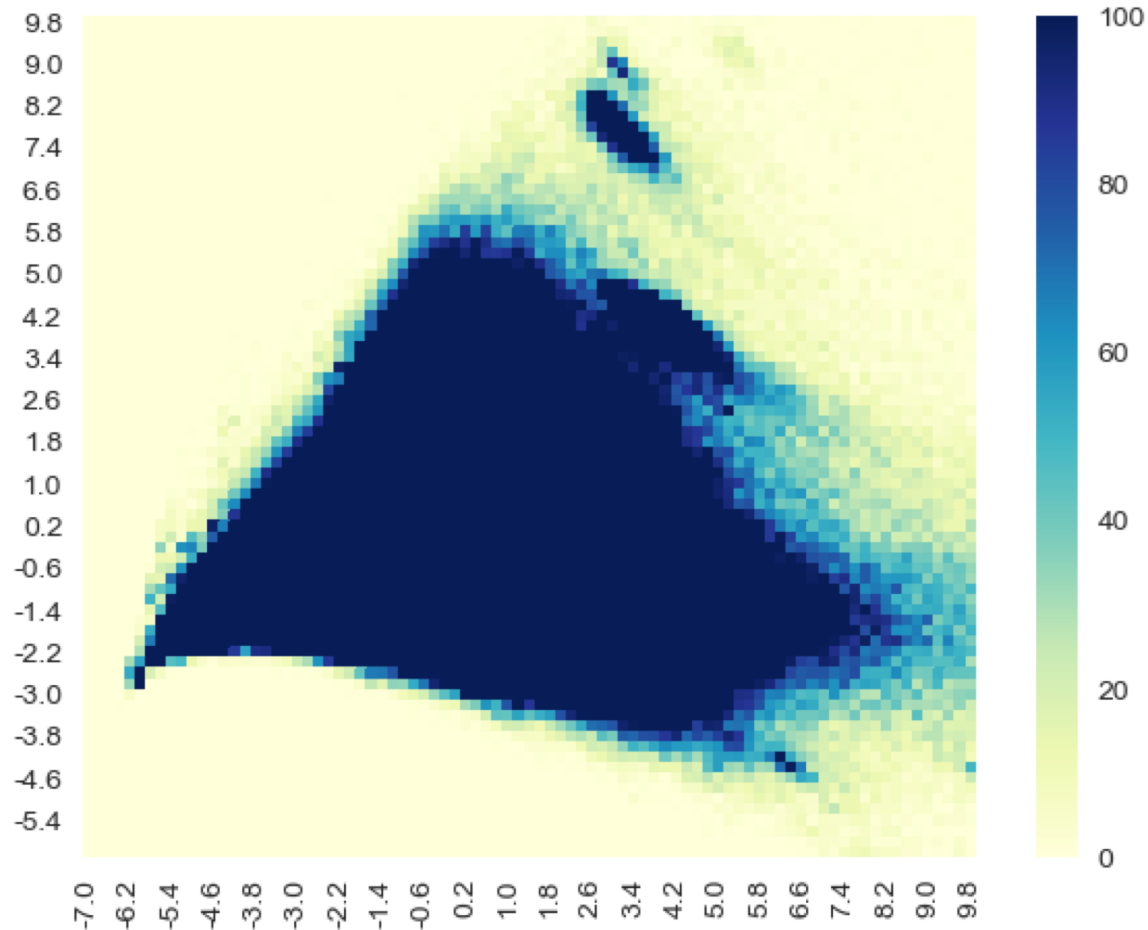
# PCA “Angelfish” Plot



# PCA “Angelfish” Plot

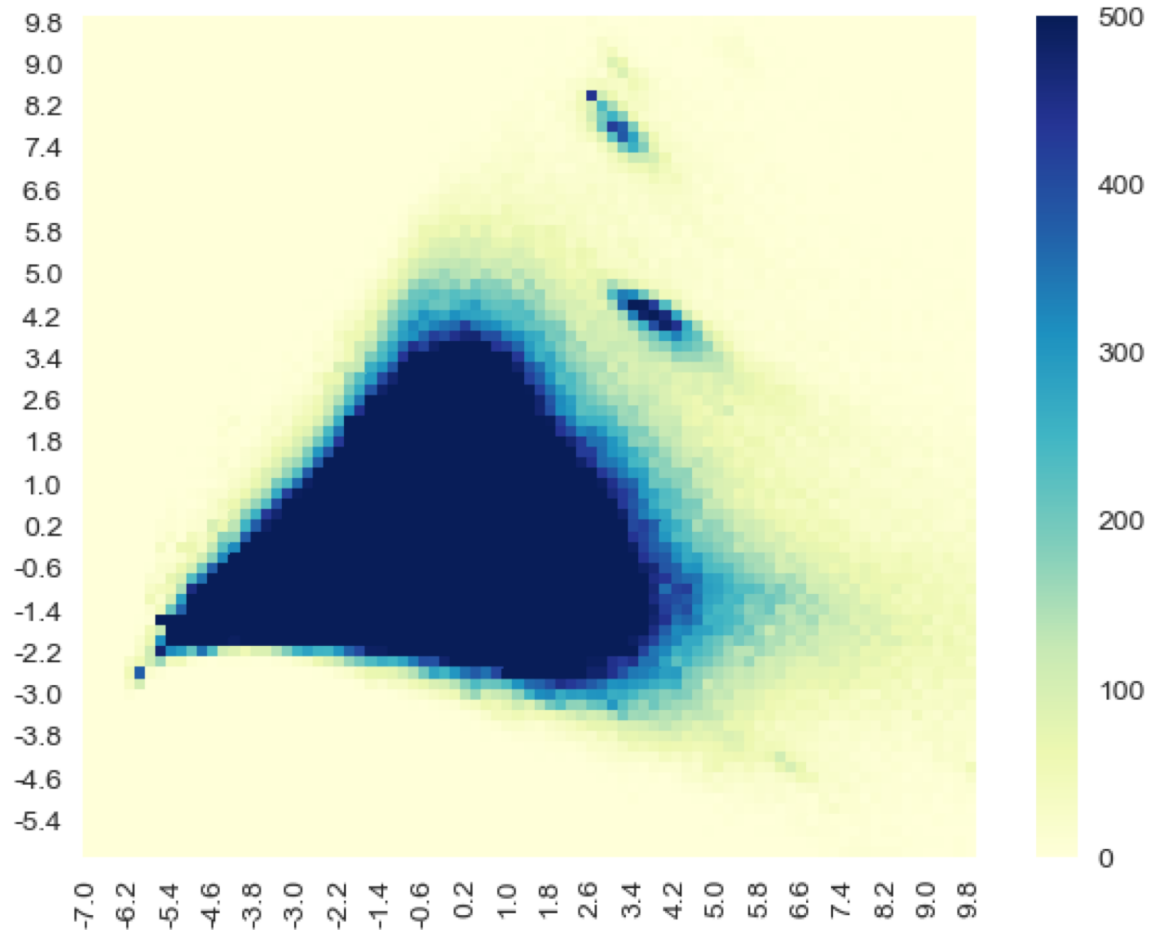


# Density estimation (dot=100pgs)

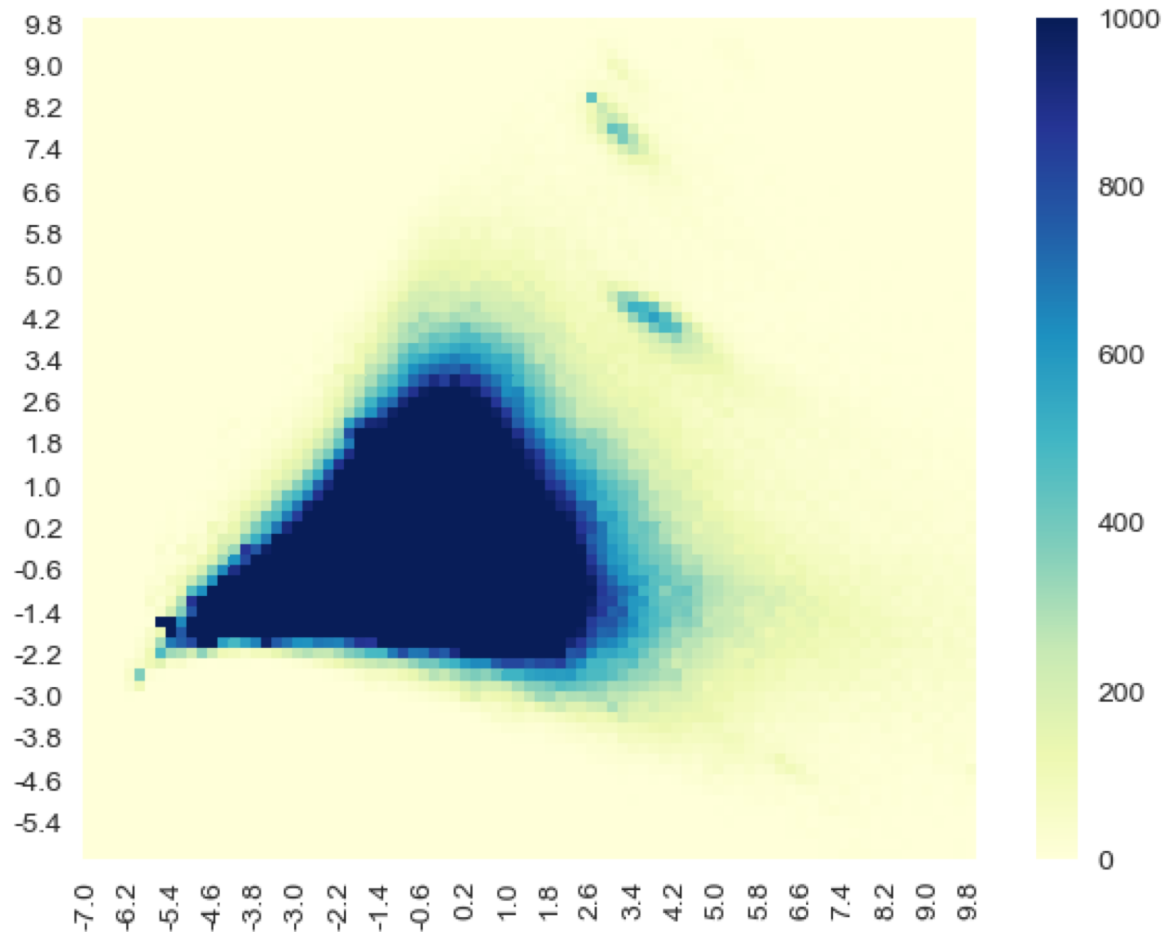




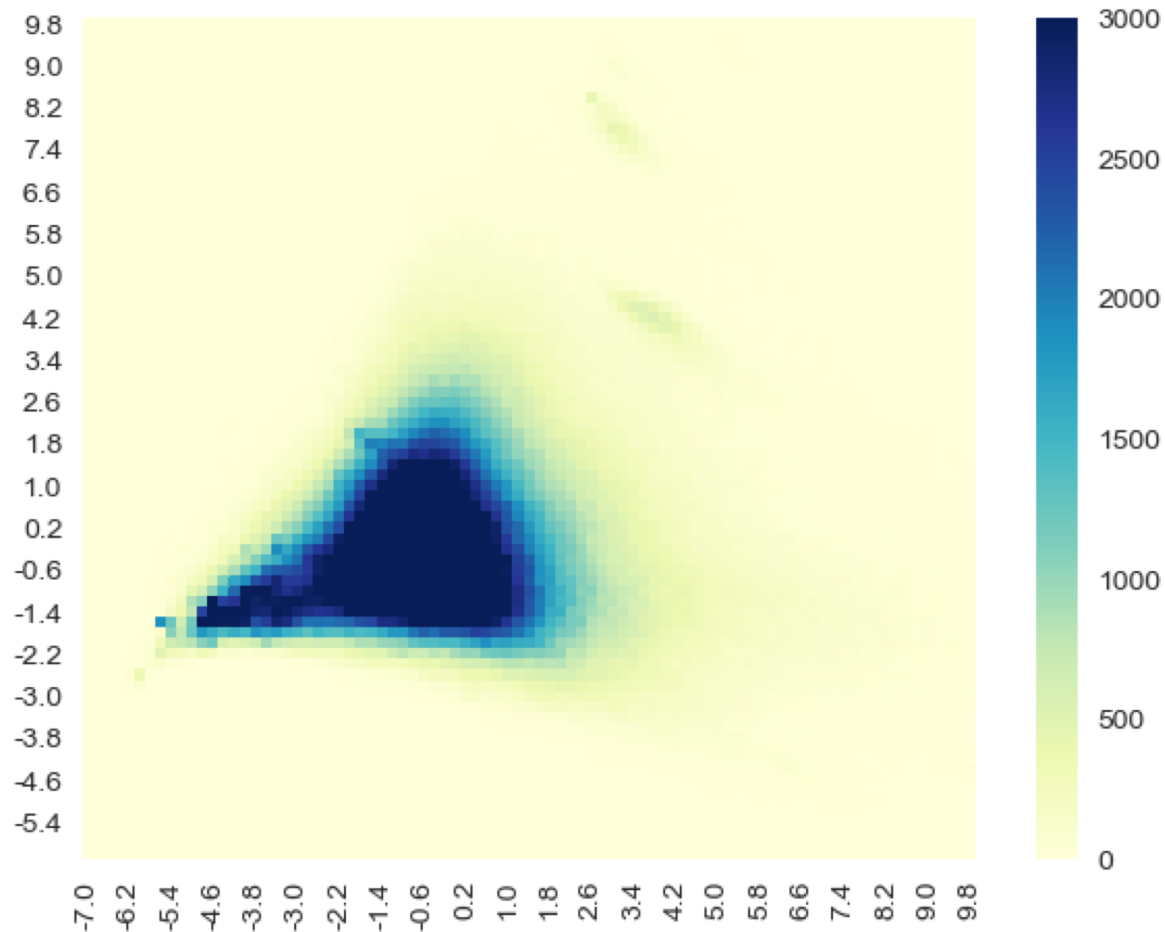
# Density estimation (dot=500pgs)



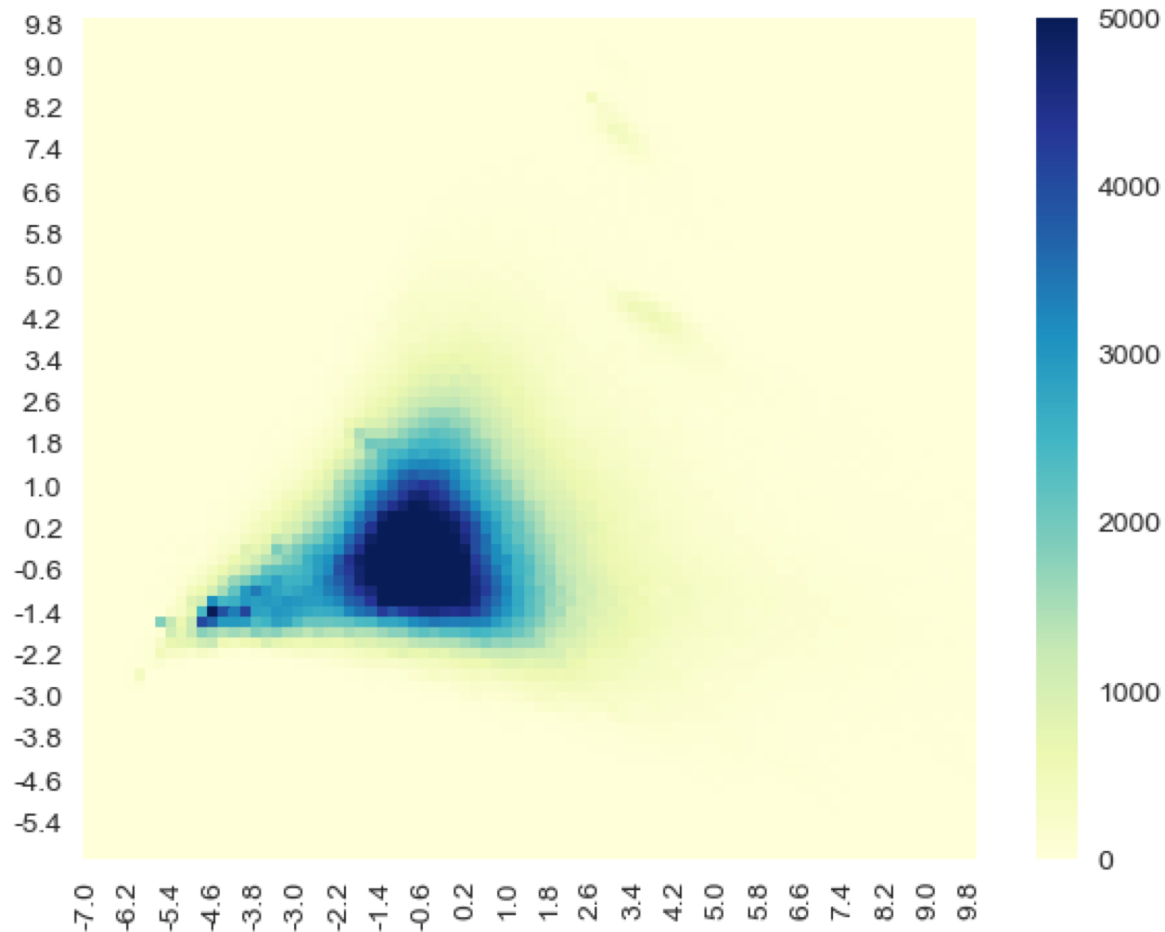
# Density estimation (dot=1 000pgs)



# Density estimation (dot=3000pgs)



# Density estimation (dot=5000pgs)



# Observations from the Density Plot

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	aa	ab	ac	ad	ae	af	ag	ah	ai	aj	ak	al	am	an	ao	ap	aq	ar	as	at	au	av	aw	ax	ay	az	ba	bb	bc	bd	be	bf	bg	bh	bi	bj	bk	bl	bm	bn	bo	bp	bq	br	bs	bt	bu	bv	bw	bx	by	bz	ca	cb	cc	cd	ce	cf	cg	ch	ci	cj	ck	cl	cm	cn	co	cp	cq	cr	cs	ct	cu	cv	cw	cx	cy	cz	da	db	dc	dd	de	df	dg	dh	di	dj	dk	dl	dm	dn	do	dp	dq	dr	ds	dt	du	dv	dw	dx	dy	dz	ea	eb	ec	ed	ee	ef	eg	eh	ei	ej	ek	el	em	en	eo	ep	eq	er	es	et	eu	ev	ew	ex	ey	ez	fa	fb	fc	fd	fe	ff	fg	fh	fi	fj	fk	fl	fm	fn	fo	fp	fq	fr	fs	ft	fu	fv	fw	fx	fy	fz	ga	gb	gc	gd	ge	gf	gg	gh	gi	gj	gk	gl	gm	gn	go	gp	gq	gr	gs	gt	gu	gv	gw	gx	gy	gz	ha	hb	hc	hd	he	hf	hg	hh	hi	hj	hk	hl	hm	hn	ho	hp	hq	hr	hs	ht	hu	hv	hw	hx	hy	hz	ia	ib	ic	id	ie	if	ig	ih	ii	ij	ik	il	im	in	io	ip	iq	ir	is	it	iu	iv	iw	ix	iy	iz	ja	jb	jc	jd	je	jf	jj	jk	jl	jm	jn	jo	jp	jq	jr	js	jt	ju	jv	jw	jx	ky	kz	la	lb	lc	ld	le	lf	lg	lh	li	lj	lk	ll	lm	ln	lo	lp	lq	lr	ls	lt	lu	lv	lw	lx	ly	lz	ma	mb	mc	md	me	mf	mg	mh	mi	mj	mk	ml	mm	mn	mo	mp	mq	mr	ms	mt	mu	mv	mw	mx	my	mz	na	nb	nc	nd	ne	nf	ng	nh	ni	nj	nk	nl	nm	nn	no	np	nq	nr	ns	nt	nu	nv	nw	nx	ny	nz	oa	ob	oc	od	oe	of	og	oh	oi	oj	ok	ol	om	on	oo	op	oq	or	os	ot	ou	ov	ow	ox	oy	oz	pa	pb	pc	pd	pe	pf	pg	ph	pi	pj	pk	pl	pm	pn	po	pp	pq	pr	ps	pt	pu	pv	pw	px	py	pz	qa	qb	qc	qd	qe	qf	qg	qh	qi	qj	qk	ql	qm	qn	qo	qp	qq	qr	qs	qt	qu	qv	qw	qx	qy	qz	ra	rb	rc	rd	re	rf	rg	rh	ri	rj	rk	rl	rm	rn	ro	rp	rq	rr	rs	rt	ru	rv	rw	rx	ry	rz	sa	sb	sc	sd	se	sf	sg	sh	si	sj	sk	sl	sm	sn	so	sp	sq	sr	ss	st	su	sv	sw	sx	sy	sz	ta	tb	tc	td	te	tf	tg	th	ti	tj	tk	tl	tm	tn	to	tp	tq	tr	ts	tt	tu	tv	tw	tx	ty	tz	ua	ub	uc	ud	ue	uf	ug	uh	ui	uj	uk	ul	um	un	uo	up	uq	ur	us	ut	uu	uv	uw	ux	uy	uz	va	vb	vc	vd	ve	vf	vg	vh	vi	vj	vk	vl	vm	vn	vo	vp	vq	vr	vs	vt	vu	vv	vw	vx	vy	vz	wa	wb	wc	wd	we	wf	wg	wh	wi	wj	wk	wl	wm	wn	wo	wp	wq	wr	ws	wt	wu	wv	ww	wx	wy	wz	xa	xb	xc	xd	xe	xf	xg	xh	xi	xj	xk	xl	xm	xn	xo	xp	xq	xr	xs	xt	xu	xv	xw	xx	xy	xz	ya	yb	yc	yd	ye	yf	yg	yh	yi	yj	yk	yl	ym	yn	yo	yp	yq	yr	ys	yt	yu	yv	yw	yx	yy	yz	za	zb	zc	zd	ze	zf	zg	zh	zi	zj	zk	zl	zm	zn	zo	zp	zq	zr	zs	zt	zu	zv	zw	zx	zy	zz
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																													

Poorly OCR'd Files

VA NPOD Consult Note

Consult Note

Missed: appt. Letter

Follow-Up

Authorization to Disclose Form

Consult?

Orthopedic Clinic Encounter

ED Legal Record

Progress Report

From Structured to "Paragraph"

VA Encounters With Questions

From Short to Structured Page Header

Urine Culture

Lab Results

PT/OT Note

Labs

Bone Density Report

Demographics

VA Note

# Topic analysis → semantic correlations

Topic	# non-relevant pages	#relevant pages
Social/family history for mental disorder	11358	975
mental status evaluation -risk of suicide	14239	122
Mental disorder symptoms and treatment history	9900	3613
Impression of mental disorder	15500	100
Lab test results [Topic 6]	11959	2406

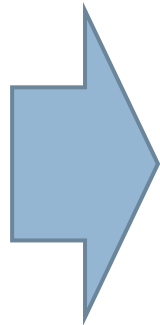
# Planting the garden: findings

47

**Content**

**Form**

**Structure**



*Technology development*

48

# Whither sublanguage analysis?



# EpiBio: Heterogeneous SSA data

49

- Geographic variation in mental health-related documentation
  - Stigma
  - Lack of details / re-coding
  
- Format → structure
  - Sectionizing
  - Semi-structured forms



# Pitt: EHR language and health equity

50

- Documentation differences for patients of different races
  - What is recorded?
  - How is it recorded?
  
- Integrating patient-generated language with clinical observations
  - Self-reported functional status
  
- Ambiguity in health language

# VA: Knowledge exchange

51

- Challenges shared by national health systems
  - ▣ Geographic and institutional variation
  - ▣ Large portion of SSA medical evidence comes from VA
  
- Cerner transition
  - ▣ Changes in documentation practice
  - ▣ Effect on NLP pipelines



# Acknowledgments

52

- Aya Zirikly (NIH)
- Guy Divita (NIH)
- Bart Desmet (NIH)
- Jona Camacho Maldonado (NIH)
- Pei-Shu Ho (NIH)
- Beth Rasch (NIH)
- Eric Fosler-Lussier (Ohio State)
- Albert Lai (WashU)



National Institutes  
of Health



Funding support from NIH  
Intramural Research Program  
and the US Social Security  
Administration.

Thank you!

[denis.griffis@nih.gov](mailto:denis.griffis@nih.gov)  
[dnewmangriffis@pitt.edu](mailto:dnewmangriffis@pitt.edu)