

A Appendix: sufficiency of l and full conditional for \mathbf{b}

Recall that the one-step-ahead conditional probability mass function in a Pólya sequence taking values in \mathbb{N} with concentration parameter α and base probability mass function Ψ is

$$p(z_i | z_{i-1}, \dots, z_1, \Psi) = \sum_{j=1}^{i-1} \frac{1}{i-1+\alpha} \mathbb{1}_{z_j}(z_i) + \frac{\alpha}{i-1+\alpha} \Psi_{z_i}. \quad (30)$$

Introducing the random variable

$$b_i \sim \text{Ber}\left(\frac{\alpha}{i-1+\alpha}\right) \quad (31)$$

we can express the one-step-ahead conditional distribution as

$$p(z_i | z_{i-1}, \dots, z_1, b_i, \Psi) = \mathbb{1}_{b_i=0} \sum_{j=1}^{i-1} \frac{1}{i-1} \mathbb{1}_{z_j}(z_i) + \mathbb{1}_{b_i=1} \Psi_{z_i}. \quad (32)$$

The joint probability mass function for $\mathbf{z} | \mathbf{b}, \Psi$ is then

$$p(\mathbf{z} | \mathbf{b}, \Psi) = \prod_{i=1}^N p(z_i | z_{i-1}, \dots, z_1, \mathbf{b}, \Psi) = \prod_{i=1}^N \left[\sum_{j=1}^{i-1} \mathbb{1}_{b_i=0} \mathbb{1}_{z_j}(z_i) + \mathbb{1}_{b_i=1} \Psi_{z_i} \right]. \quad (33)$$

Note that $\mathbb{1}_{b_i=0} = 1 \iff \mathbb{1}_{b_i=1} = 0$ and vice versa. Thus each term in the product for $\mathbf{z} | \mathbf{b}, \Psi$ only has one component, and we may express $\mathbf{z} | \mathbf{b}, \Psi$ as

$$p(\mathbf{z} | \mathbf{b}, \Psi) = \underbrace{\prod_{\substack{i=1 \\ b_i \neq 1}}^N \sum_{j=1}^{i-1} \frac{1}{i-1} \mathbb{1}_{z_j}(z_i)}_{\text{doesn't enter posterior}} \prod_{\substack{i=1 \\ b_i=1}}^N \prod_{k=1}^{\infty} \Psi_k^{\mathbb{1}_k(z_i)} \quad (34)$$

where we have re-expressed the probability mass function of Ψ in a form that emphasizes conjugacy. Thus for any prior, the posterior will only depend on the likelihood of the values of z_i for which $b_i = 1$. The sufficient statistic is

$$l_k = \sum_{\substack{i=1 \\ b_i=1}}^N \mathbb{1}_{z_i=k}. \quad (35)$$

Next, for a given $i' \in \{1, \dots, N\}$, we can calculate the posterior of a component $b_{i'}$ as

$$\mathbb{P}(b_{i'} = 1 | \mathbf{z}, \Psi, \mathbf{b}_{-i'}) \propto \left(\frac{\alpha}{i'-1+\alpha} \right) \prod_{\substack{i=1 \\ b_i \neq 1}}^N \sum_{j=1}^{i-1} \frac{1}{i-1} \mathbb{1}_{z_j}(z_i) \prod_{\substack{i=1 \\ b_i=1}}^N \Psi_{z_i} \quad (36)$$

$$\propto \alpha \Psi_{z_{i'}} \quad (37)$$

$$\mathbb{P}(b_{i'} = 0 | \mathbf{z}, \Psi, \mathbf{b}_{-i'}) \propto \left(\frac{i'-1}{i'-1+\alpha} \right) \prod_{\substack{i=1 \\ b_i \neq 1}}^N \sum_{j=1}^{i-1} \frac{1}{i-1} \mathbb{1}_{z_j}(z_i) \prod_{\substack{i=1 \\ b_i=1}}^N \Psi_{z_i} \quad (38)$$

$$\propto \sum_{i=1}^{i'-1} \mathbb{1}_{z_i}(z_{i'}) \quad (39)$$

where we have divided both expressions by

$$\frac{1}{i'-1+\alpha} \prod_{\substack{i=1 \\ b_i \neq 1 \\ i \neq i'}}^N \sum_{j=1}^{i-1} \frac{1}{i-1} \mathbb{1}_{z_j}(z_i) \prod_{\substack{i=1 \\ b_i=1 \\ i \neq i'}}^N \Psi_{z_i} \quad (40)$$

which is constant with respect to $b_{i'}$. Note that full conditionally, we have $b_i \perp\!\!\!\perp b_{i'}$ for $i \neq i'$. This gives the desired expressions and concludes the derivation.

B Appendix: full conditional for Ψ

Before proceeding with the derivation, we first comment on Proposition 1 and differences between the GEM distribution and Dirichlet process, which otherwise appear superficially similar. The GEM distribution $\Psi^{\text{GEM}} \sim \text{GEM}(\gamma)$ is defined as

$$\Psi_k^{\text{GEM}} = \varsigma_k \prod_{i=1}^{k-1} (1 - \varsigma_i) \quad \varsigma_k^{\text{GEM}} \sim \text{Beta}(1, \gamma). \quad (41)$$

On the other hand, a Dirichlet process $\Psi^{\text{DP}} \sim \text{DP}(\gamma, F)$ is defined as

$$\Psi^{\text{DP}} = \sum_{k=1}^{\infty} \pi_k \delta_{\vartheta_k} \quad \vartheta_k \sim F \quad \pi_k = \varsigma_k \prod_{i=1}^{k-1} (1 - \varsigma_i) \quad \varsigma_k \sim \text{Beta}(1, \gamma). \quad (42)$$

From a Bayesian perspective, this extra stage—the presence of ϑ_k —prevents one from applying standard results on conjugacy of Dirichlet processes. The joint distribution of a finite set of states $(\Psi_{k_1}^{\text{GEM}}, \dots, \Psi_{k_K}^{\text{GEM}})$ does not admit a closed-form expression, so we seek to derive the posterior conditional in a different way.

Rather than proving conjugacy for $(\Psi_{k_1}^{\text{GEM}}, \dots, \Psi_{k_K}^{\text{GEM}})$ directly, we look for a larger finite-dimensional distribution within which $(\Psi_{k_1}^{\text{GEM}}, \dots, \Psi_{k_K}^{\text{GEM}})$ sits that has better conjugacy properties. The *generalized Dirichlet* distribution of [Connor and Mosimann \(1969\)](#) fulfills this criteria. The conjugacy relationship we seek follows from the general property that conditioning and marginalization commute. This will be shown to yield the posterior

$$\Psi_k^{\text{GEM}} = \varsigma_k \prod_{i=1}^{k-1} (1 - \varsigma_i) \quad \varsigma_k \sim \text{Beta}(a_k^{(\Psi)}, b_k^{(\Psi)}) \quad a_k^{(\Psi)} = 1 + l_k \quad b_k^{(\Psi)} = \gamma + \sum_{i=k+1}^{\infty} l_i. \quad (43)$$

For comparison, a posterior Dirichlet process is given by

$$\Psi^{\text{DP}} = \sum_{k=1}^{\infty} \pi_k \delta_{\vartheta_k} \quad \vartheta_k \sim \frac{n}{\alpha + n} \mathbf{l} + \frac{\alpha}{\alpha + n} F \quad \pi_k = \varsigma_k \prod_{i=1}^{k-1} (1 - \varsigma_i) \quad \varsigma_k \sim \text{Beta}(1, \gamma + n) \quad (44)$$

which shows that this relatively mild difference in the prior yields a posterior of a rather different form.

We now proceed to formally calculate this posterior distribution, starting from a GEM prior and discrete likelihood. Since we are working in a nonparametric setting, we begin by introducing the necessary formalism. We then introduce our finite-dimensional approximating prior and compute the posterior under it. For this, we use commutativity of conditioning and marginalization to deduce the full infinite-dimensional posterior.

Definition 2 (Preliminaries). Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. Let $\mathcal{M}_s(\mathbb{N})$ be the space of signed measures, equipped with the topology of weak convergence. Let $\mathcal{M}_1(\mathbb{N}) \subset \mathcal{M}_s(\mathbb{N})$ be the space of probability measures over \mathbb{N} , and identify $\mathcal{M}_1(\mathbb{N})$ with the probability simplex by the homeomorphism $\mathcal{M}_1(\mathbb{N}) \cong \{\mathbf{x} \in \ell^1 : \forall i, x_i > 0, \sum_{i=1}^{\infty} x_i = 1\}$. Let $N \in \mathbb{N}$, let $\mathbf{x} \in \mathbb{N}^N$, and let $\mathbf{l} \in \mathbb{N}^N$ be its empirical counts, defined by $\mathbf{l} = \sum_{i=1}^N \mathbf{1}_{x_i}$ where $\mathbf{1}_{x_i}$ is equal to 1 for coordinate x_i and 0 for all other coordinate. Let $\gamma \in \mathbb{R}^+$. Recall that \mathbb{N}^N and $\mathcal{M}_1(\mathbb{N})$, endowed with the discrete topology and topology of weak convergence, respectively, are both Polish spaces—hence, the Disintegration Theorem ([Ambrosio et al. \(2005\)](#), Theorem 5.3.1; [Bogachev \(2007\)](#), Corollary 10.4.15) holds in both spaces. We associate each random variable $y : \Omega \rightarrow Y$ with its pushforward probability measure $\pi_y(A_y) = [y_* \mathbb{P}](A_y) = \mathbb{P}[y^{-1}(A_y)]$, and each conditional random variables $\theta | y : \Omega \times Y \rightarrow \Theta$ with its pushforward regular conditional probability measure $\pi_{y|\theta}(A_y | \theta) = [(y | \theta)_* \mathbb{P}](A_y) = \mathbb{P}[(y | \theta)^{-1}(A_y)]$, where the preimage is taken with respect to y .

Definition 3 (Discrete likelihood). For all $\Psi \in \mathcal{M}_1(\mathbb{N})$, define the conditional random variable $\mathbf{x} \mid \Psi : \Omega \times \mathcal{M}_1(\mathbb{N}) \rightarrow \mathbb{N}^N$ by its probability mass function

$$p(\mathbf{x} \mid \Psi) = \prod_{i=1}^N \prod_{k=1}^{\infty} \Psi_k^{\mathbb{1}_k(x_i)}. \quad (45)$$

We say $\mathbf{x} \mid \Psi \sim \text{Discrete}(\Psi)$.

Definition 4 (GEM). Let $\Psi : \Omega \rightarrow \mathcal{M}_1(\mathbb{N})$ be a random variable defined by

$$\Psi_k = \varsigma_k \prod_{i=1}^{k-1} (1 - \varsigma_i) \quad \varsigma_k \sim \text{Beta}(1, \gamma). \quad (46)$$

We say $\Psi \sim \text{GEM}(\gamma)$.

Definition 5 (Finite GEM). Let $\Psi : \Omega \rightarrow \mathcal{M}_1(\mathbb{N})$ be a random variable defined by

$$\Psi_k = \varsigma_k \prod_{i=1}^{k-1} (1 - \varsigma_i) \quad \varsigma_k \sim \text{Beta}(1, \gamma) \quad \varsigma_K = 1. \quad (47)$$

We say $\Psi \sim \text{FGEM}(\gamma, K)$.

Definition 6 (Posterior). Let $\Psi \mid \mathbf{x}$ be the unique conditional random variable given by the Disintegration Theorem, where uniqueness follows from almost sure uniqueness by virtue of the marginal measure $\pi_{\mathbf{x}}(\cdot) = \int_{\mathcal{M}_1(\mathbb{N})} \pi_{\mathbf{x} \mid \Psi}(\cdot \mid \Psi) d\pi_{\Psi}$ being absolutely continuous with respect to the counting measure on \mathbb{N}^N , which has no non-empty null sets.

Result 7. Let $\mathbf{x} \mid \Psi \sim \text{Discrete}(\Psi)$. Let $\mathbf{x} \in \mathbb{N}^N$, and let $K > \sup \mathbf{x}$. Let $\Psi \sim \text{FGEM}(\gamma, K)$. Then for any \mathbf{x} with empirical counts \mathbf{l} , we have that $\Psi \mid \mathbf{x} : \Omega \times \mathbb{N}^N \rightarrow \mathcal{M}_1(\mathbb{N})$ is a conditional random variable defined by

$$\Psi_k = \varsigma_k \prod_{i=1}^{k-1} (1 - \varsigma_i) \quad \varsigma_k \sim \text{Beta}(a_k^{(\Psi)}, b_k^{(\Psi)}) \quad \varsigma_K = 1 \quad (48)$$

where

$$a_k^{(\Psi)} = 1 + l_k \quad b_k^{(\Psi)} = \gamma + \sum_{i=k+1}^K l_i. \quad (49)$$

Proof. It is shown by Connor and Mosimann (1969) that $\Psi \sim \text{FGEM}(\gamma, K)$ is a special case of the *generalized Dirichlet* distribution, which admits a general stick-breaking representation. Thus, its probability density function is

$$f(\Psi) \propto \Psi_K^{\gamma-1} \prod_{k=1}^{K-1} \left[\sum_{k'=1}^K \Psi_{k'} \right]^{-1} \quad (50)$$

which we have expressed in a simplified form. By conjugacy, for a given \mathbf{x} and associated \mathbf{l} the posterior probability density is

$$f(\Psi \mid \mathbf{x}) \propto \Psi_K^{(\gamma+l_K)-1} \prod_{k=1}^{K-1} \left[\Psi_k^{(1+l_k)-1} \left[\sum_{k'=k}^K \Psi_{k'} \right]^{\gamma + \sum_{i=k}^K l_i - [(1+l_k) + \gamma + \sum_{i=k+1}^K l_i]} \right] \quad (51)$$

which is again a generalized Dirichlet admitting the necessary stick-breaking representation, which we have expressed in a form that emphasizes its posterior hyperparameters. \square

Remark 8. It is now clear that the assumption $\boldsymbol{x} \mid \Psi \sim \text{Discrete}(\Psi)$ is indeed taken without loss of generality, because if we instead took $\boldsymbol{x} \mid \Psi$ to be given by a Pólya sequence, then by sufficiency the prior-to-posterior map would be identical.

Proposition 1. Without loss of generality, suppose

$$\Psi \sim \text{GEM}(\gamma) \qquad \boldsymbol{x} \mid \Psi \sim \text{Discrete}(\Psi). \qquad (52)$$

Then $\Psi \mid \boldsymbol{x}$ is given by

$$\Psi_k = \varsigma_k \prod_{i=1}^{k-1} (1 - \varsigma_i) \quad \varsigma_k \sim \text{Beta}(a_k^{(\Psi)}, b_k^{(\Psi)}) \quad a_k^{(\Psi)} = 1 + l_k \quad b_k^{(\Psi)} = \gamma + \sum_{i=k+1}^{\infty} l_i \qquad (53)$$

where l are the empirical counts of \boldsymbol{x} .

Proof. Let $I \subset \mathbb{N}$ be an arbitrary finite index set, and let $\Psi_I \mid \boldsymbol{x}$ be the finite-dimensional marginal projection of $\Psi \mid \boldsymbol{x}$ onto the coordinates contained in I . Let $K > \sup I$, let $\Psi^{(K)} \mid \boldsymbol{x}$ be the posterior conditional random variable under $\Psi^{(K)} \sim \text{FGEM}(\gamma, K)$, and let $\Psi_I^{(K)} \mid \boldsymbol{x}$ be the marginal consisting of those coordinates contained in I . By construction, $\Psi_I^{(K)} \mid \boldsymbol{x}$ equals $\Psi_I \mid \boldsymbol{x}$ in distribution. Since by the Disintegration Theorem, conditioning and marginalization commute, the set I is arbitrary, and $\Psi \mid \boldsymbol{x}$ is uniquely determined by its finite-dimensional marginal projections, the claim follows. \square

C Appendix: quantile summary of topics for AP

Here we display a multi-quantile summary for AP, obtained by ranking all topics with at least 100 tokens by their total number of tokens, computing the $\varpi = 100\%$, 75%, 50%, 25%, and 5% quantiles. We compute the five topics closest to each quantile by number of tokens, and display their top-eight words.

AP partially collapsed $\varpi = 100\%$	k	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
	$n_{k,\cdot}$	93 207	57 249	15 874	13 360	10 176
		week	people	police	percent	trial
		made	years	people	year	court
		president	year	killed	prices	charges
		officials	time	man	economic	case
		tuesday	don	officials	economy	judge
		million	back	city	rate	attorney
	thursday	day	shot	increase	prison	
	national	home	authorities	report	jury	
AP partially collapsed $\varpi = 75\%$	k	Topic 54	Topic 55	Topic 56	Topic 57	Topic 58
	$n_{k,\cdot}$	1 055	1 032	1 025	1 014	1 013
		children	north	hostages	aids	percent
		parents	walsh	red	virus	poll
		child	reagan	release	blood	survey
		ms	iran	held	disease	points
		year	contra	hostage	drug	found
		mother	documents	anderson	infected	surveys
	boys	gesell	gunmen	immune	margin	
	girl	arms	thursday	health	reported	
AP partially collapsed $\varpi = 50\%$	k	Topic 108	Topic 109	Topic 110	Topic 111	Topic 112
	$n_{k,\cdot}$	473	472	451	446	436
		abortion	women	solidarity	waste	train
		souter	club	walesa	garbage	railroad
		anti	members	poland	recycling	cars
		state	men	polish	city	trains
		women	male	government	ash	transportatio
		abortions	membership	mazowiecki	trash	skinner
	rights	female	jaruzelski	state	transit	
	hampshire	black	talks	dump	policy	
AP partially collapsed $\varpi = 25\%$	k	Topic 162	Topic 163	Topic 164	Topic 165	Topic 166
	$n_{k,\cdot}$	193	187	185	184	184
		health	wine	miners	dixon	barry
		care	warmus	mine	yates	moore
		spe	solomon	coal	count	jackson
		bc	california	mines	tosh	mayor
		weight	bar	hull	rogers	statehood
		american	gallo	pittston	rig	mr
	diet	test	benefits	russell	gregory	
	cholesterol	questions	platform	cookies	room	
AP partially collapsed $\varpi = 5\%$	k	Topic 206	Topic 207	Topic 208	Topic 209	Topic 210
	$n_{k,\cdot}$	117	115	112	111	111
		pageant	mall	roberts	stuart	gold
		miss	malls	shell	lawn	polaroid
		cereal	pinochet	boigny	dea	shamrock
		boxes	shopping	houphouet	boston	fields
		contestants	downtown	travelers	ruth	consolidated
		box	park	leonard	yankees	suit
	america	oak	arsenal	foundation	proposals	
	bruce	usa	oil	richman	mining	

AP	k $n_{k,\bullet}$	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5		
		90 497	18 626	10 832	9 923	9 430		
AP	direct assignment $\varpi = 100\%$	year	years	police	dollar	percent		
		people	year	people	market	year		
		time	people	killed	stock	rose		
		president	time	government	yen	sales		
		years	don	reported	index	million		
		made	home	today	late	billion		
		state	day	capital	trading	month		
		week	back	violence	exchange	reported		
		AP	k $n_{k,\bullet}$	Topic 93	Topic 94	Topic 95	Topic 96	Topic 97
				784	757	753	745	738
AP	direct assignment $\varpi = 75\%$	keating	bus	eastern	united	smoking		
		deconcini	driver	pilots	states	cigarettes		
		lincoln	train	airline	nations	farmers		
		senators	greyhound	orion	resolution	tobacco		
		regulators	accident	air	international	ban		
		meeting	passengers	union	plo	insurance		
		committee	railroad	airlines	mission	batus		
		gray	passenger	service	assembly	smokers		
		AP	k $n_{k,\bullet}$	Topic 186	Topic 187	Topic 188	Topic 189	Topic 190
				346	338	338	338	334
AP	direct assignment $\varpi = 50\%$	power	cable	conservatives	water	dental		
		franc	television	flag	dam	funds		
		jersey	nbc	conservative	river	claims		
		bradley	tempo	amendment	area	plough		
		utility	hsn	speaker	reservoir	oral		
		wppss	industry	darman	savannah	counter		
		utilities	subscribers	kemp	corps	embassy		
		west	tv	republicans	canyon	mid		
		AP	k $n_{k,\bullet}$	Topic 279	Topic 280	Topic 281	Topic 282	Topic 283
				220	219	219	219	218
AP	direct assignment $\varpi = 25\%$	fernandez	water	bloom	canadian	election		
		fdic	lake	minnick	lee	grenada		
		republicbank	mussels	walters	ritalin	boigny		
		weicker	neill	lawyer	murphy	houphouet		
		virginia	erie	athletes	domestic	gairy		
		ruth	problem	college	security	coast		
		robinson	plant	suspect	woods	nov		
		station	north	signing	radio	failed		
		AP	k $n_{k,\bullet}$	Topic 354	Topic 355	Topic 356	Topic 357	Topic 358
				133	133	133	132	132
AP	direct assignment $\varpi = 5\%$	machine	young	count	reynolds	turkey		
		stop	johnston	forman	premier	department		
		reed	golf	festival	bond	bird		
		gun	notes	rig	release	cooking		
		chief	bodies	arts	news	wash		
		sununu	homes	hughes	regulated	bacteria		
		geneva	call	lights	address	stuffed		
		formal	shortage	staged	petition	adams		

D Appendix: quantile summary of topics for CGCBIB

Here we display a multi-quantile summary for CGCBIB, obtained by ranking all topics with at least 100 tokens by their total number of tokens, computing the $\varpi = 100\%$, 75% , 50% , 25% , and 5% quantiles. We compute the five topics closest to each quantile by number of tokens, and display their top-eight words.

CGCBIB partially collapsed $\varpi = 100\%$	k	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
	$n_{k,\cdot}$	110 702	58 811	27 084	21 215	19 832
		elegans	elegans	elegans	gene	mutations
		caenorhabditi	protein	genetic	elegans	gene
		nematode	caenorhabditi	molecular	sequence	mutants
		results	gene	development	protein	genes
		found	function	caenorhabditi	caenorhabditi	mutant
		show	proteins	nematode	amino	elegans
	observed	required	studies	cdna	caenorhabditi	
	specific	show	model	acid	alleles	
CGCBIB partially collapsed $\varpi = 75\%$	k	Topic 54	Topic 55	Topic 56	Topic 57	Topic 58
	$n_{k,\cdot}$	2 166	2 061	2 048	2 040	2 025
		germ	egl	emb	spe	wnt
		germline	egg	temperature	sperm	mom
		early	laying	mutants	fer	pop
		granules	serotonin	sensitive	spermatozoa	signaling
		cells	neurons	zyg	membrane	bar
		embryos	cat	maternal	spermatids	pathway
	somatic	dopamine	expression	spermatogenes	lin	
	line	mutants	embryonic	pseudopod	wrm	
CGCBIB partially collapsed $\varpi = 50\%$	k	Topic 109	Topic 110	Topic 111	Topic 112	Topic 113
	$n_{k,\cdot}$	930	916	915	900	893
		vit	binding	kinesin	growth	eat
		yolk	affinity	klp	survival	pharyngeal
		vitellogenin	site	transport	mortality	pharynx
		genes	activity	motor	population	pumping
		yp	sites	ift	rate	inx
		proteins	avermectin	cilia	populations	gap
	vpe	elegans	dynein	parameter	feeding	
	lrp	membrane	movement	size	junctions	
CGCBIB partially collapsed $\varpi = 25\%$	k	Topic 164	Topic 165	Topic 166	Topic 167	Topic 168
	$n_{k,\cdot}$	386	369	368	364	360
		mlc	dom	innate	vha	ife
		mel	effects	immune	atpase	cap
		myosin	humic	immunity	subunit	eife
		nmy	pyrene	abf	genes	capping
		chain	effect	lys	vacuolar	cel
		elongation	bioconcentrat	toll	subunits	gtp
	rho	dissolved	antimicrobial	atpases	isoforms	
	phosphatase	substances	pathway	type	rna	
CGCBIB partially collapsed $\varpi = 5\%$	k	Topic 208	Topic 209	Topic 210	Topic 211	Topic 212
	$n_{k,\cdot}$	141	141	140	140	136
		ubq	asp	da	ion	hcf
		gc	salmonella	cl	diet	cehcf
		tbp	poona	fli	relative	vp
		footprints	enterica	gs	xpa	ldb
		oscillin	clp	db	groups	cell
		tlf	serotype	glu	carbon	mammalian
	ubiquitin	necrotic	phospholipid	characteristi	phosphorylati	
	tata	mug	tg	atoms	neural	

CGCBIB	k $n_{k,\bullet}$	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
		65 059	41 005	33 714	27 221	22 813
CGCBIB	direct assignment $\varpi = 100\%$	elegans caenorhabditi protein gene function proteins required show	elegans genetic caenorhabditi nematode molecular development studies model	elegans caenorhabditi nematode results observed high type effect	mutations elegans gene mutants genes caenorhabditi mutant function	elegans gene sequence caenorhabditi protein amino cdna acid
		Topic 68	Topic 69	Topic 70	Topic 71	Topic 72
CGCBIB	k $n_{k,\bullet}$	1 921	1 894	1 836	1 828	1 776
		loci genetic strains lines life mutations mutation inbred	worm elegans research caenorhabditi brenner years nematode biology	cell epithelial junctions membrane cells dlg hmp exc	alpha gpa egl signaling protein goa rgs proteins	unc gaba receptors receptor resistance lev levamisole cholinergic
CGCBIB	k $n_{k,\bullet}$	Topic 137	Topic 138	Topic 139	Topic 140	Topic 141
		782	779	779	763	756
CGCBIB	direct assignment $\varpi = 50\%$	cell dimensional microscopy embryo analysis system computer time	hsp heat shock chaperone small proteins crystallin hsps	survival mortality model data gompertz parameter population rate	acid amino acids nematode glycine briggsae cytochrome multiple	yeast cerevisiae saccharomyces pombe cell budding schizosacchar cycle
		Topic 206	Topic 207	Topic 208	Topic 209	Topic 210
CGCBIB	k $n_{k,\bullet}$	364	359	358	343	343
		activity lh activities juvenile nematodes antiallatal hormone insect	pgp mrp aat cells glycoprotein mammalian resistance glycoproteins	telomere telomeres ceh yeast nematode mrt telomerase telomeric	gcy guanylyl cyclase wee ase receptor cyclases gfp	mediator med sop transcription development pvl transcription dhp
CGCBIB	k $n_{k,\bullet}$	Topic 261	Topic 262	Topic 263	Topic 264	Topic 265
		164	164	159	157	156
CGCBIB	direct assignment $\varpi = 5\%$	cog wd repeat connection native response worm nr	atp structures oligomerizati family binding members stability mechanism	calcineurin egg bovine laying hg white haemin phosphatase	selection flow separation redundancy flows directional solution period	srl rol threshold ra energy free external experimental

E Appendix: quantile summary of topics for NEURIPS

Here we display a multi-quantile summary for NeurIPS, obtained by ranking all topics with at least 100 tokens by their total number of tokens, computing the $\varpi = 100\%$, 75%, 50%, 25%, and 5% quantiles. We compute the five topics closest to each quantile by number of tokens, and display their top-eight words.

NeurIPS partially collapsed $\varpi = 100\%$	k	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
	$n_{k,\cdot}$	182 743	162 355	129 745	52 356	44 155
		system	function	number	model	training
		information	case	result	neural	set
		approach	result	small	result	data
		set	term	values	system	test
		problem	parameter	order	activity	performance
		research	neural	large	input	number
	computer	form	effect	pattern	result	
	single	defined	high	function	error	
NeurIPS partially collapsed $\varpi = 75\%$	k	Topic 148	Topic 149	Topic 150	Topic 151	Topic 152
	$n_{k,\cdot}$	2 585	2 585	2 574	2 559	2 549
		genetic	delay	bengio	fig	matching
		algorithm	bifurcation	output	properties	model
		population	oscillation	dependencies	proc	point
		fitness	point	input	step	correspondenc
		string	stability	experiment	range	match
		generation	fixed	frasconi	structure	problem
	bit	limit	term	calculation	set	
	function	hopf	information	illinois	object	
NeurIPS partially collapsed $\varpi = 50\%$	k	Topic 297	Topic 298	Topic 299	Topic 300	Topic 301
	$n_{k,\cdot}$	1 310	1 309	1 309	1 297	1 295
		vor	routing	speaker	delay	memory
		storage	load	recognition	input	action
		anastasio	network	normalization	transition	states
		responses	path	male	window	agent
		velocity	packet	feature	width	sensing
		pan	traffic	female	connection	loop
	rotation	shortest	mntn	information	history	
	vestibular	policy	ntn	temporal	mdp	
NeurIPS partially collapsed $\varpi = 25\%$	k	Topic 446	Topic 447	Topic 448	Topic 449	Topic 450
	$n_{k,\cdot}$	748	748	746	739	735
		composite	psom	limited	tau	cmm
		mdp	robot	interconnect	hypothesis	speed
		action	camera	fan	mansour	particle
		elemental	set	shunting	growth	particles
		optimal	pointing	modularity	coefficient	pattern
		payoff	coordinates	collective	function	presence
	solution	basis	linear	stem	method	
	mdt	ritter	unit	large	card	
NeurIPS partially collapsed $\varpi = 5\%$	k	Topic 566	Topic 567	Topic 568	Topic 569	Topic 570
	$n_{k,\cdot}$	396	385	383	379	372
		morph	minimal	visualization	periodic	machine
		kernel	root	high	period	capacity
		parent	biases	low	coefficient	path
		human	attribute	diagram	primitive	trouble
		busey	remove	visualizing	homogeneous	high
		similar	rumelhart	graphic	tst	task
	exemplar	row	fund	mhaskar	increasing	
	distinctivene	exponential	window	chain	measures	

NeurIPS subcluster split-merge $\varpi = 100\%$	k $n_{k,\bullet}$	Topic 6 473 770	Topic 2 93 435	Topic 1 52 418	Topic 13 50 965	Topic 62 41 565
		network model learning function input neural algorithm set	network unit input learning training weight neural output	model neuron input network cell system unit visual	model data parameter network algorithm mixture function gaussian	function network bound dimension learning result number set
NeurIPS subcluster split-merge $\varpi = 75\%$	k $n_{k,\bullet}$	Topic 440 2 678	Topic 170 2 657	Topic 334 2 643	Topic 418 2 636	Topic 312 2 622
		learning critic function actor algorithm system control model	movement visual vector image model location eye map	motion unit direction model stage input network cell	learning algorithm advantage system function policy control	cell correlation neuron model unit interaction firing set
NeurIPS subcluster split-merge $\varpi = 50\%$	k $n_{k,\bullet}$	Topic 378 1 032	Topic 322 1 028	Topic 82 1 013	Topic 344 1 009	Topic 414 1 006
		iii cell network neural response model point fixed	cell spike unit function firing result transfer sorting	model response neural escape interneuron cockroach leg input	form word phone input network system training meaning	component algorithm sources analysis data noise orientation spatial
NeurIPS subcluster split-merge $\varpi = 25\%$	k $n_{k,\bullet}$	Topic 220 728	Topic 341 723	Topic 441 723	Topic 447 722	Topic 308 721
		aspect object view node learning network weight equation	element pairing grouping group saliency contour computation optimal	network neural constraint match learn problem initial row	input unit spike layer learning model predict prediction	traffic waiting elevator appeared application compared department found
NeurIPS subcluster split-merge $\varpi = 5\%$	k $n_{k,\bullet}$	Topic 259 509	Topic 246 507	Topic 195 506	Topic 245 503	Topic 293 503
		input output activation data encoded function hidden model	network neural task link food nodes output recurrent	network symbol vtp learning phrases sentences vpp classificatio	network equation neuron moment neural approximation ohira stochastic	network function adaptation algorithm prediction projection neural training

F Appendix: topics produced by Algorithm 2 on PUBMED

Here we show top eight words for each topic together with total number of tokens assigned, which is shown at the top of each table. We display all topics containing at least eight unique word tokens.

k	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
$n_{k,\bullet}$	47 322 709	40 229 486	34 685 122	30 795 166	30 707 144
PubMed	care	age	model	cell	gene
	health	risk	data	expression	protein
	patient	children	system	growth	dna
	medical	year	time	protein	expression
	research	women	analysis	factor	sequence
	clinical	patient	effect	receptor	genes
	system	factor	test	kinase	rna
	cost	population	field	beta	region
k	Topic 6	Topic 7	Topic 8	Topic 9	Topic 10
$n_{k,\bullet}$	28 510 997	27 277 306	26 709 116	26 408 263	25 200 662
PubMed	cell	cancer	patient	rat	cell
	il	tumor	treatment	receptor	electron
	cd	patient	mg	effect	muscle
	mice	carcinoma	drug	neuron	tissue
	antigen	cell	effect	brain	fiber
	human	breast	therapy	activity	rat
	lymphocytes	survival	dose	stimulation	development
	immune	tumour	day	induced	microscopy
k	Topic 11	Topic 12	Topic 13	Topic 14	Topic 15
$n_{k,\bullet}$	24 856 624	24 750 437	24 607 618	24 482 090	22 956 810
PubMed	patient	patient	blood	patient	infection
	surgery	artery	pressure	disease	virus
	complication	heart	flow	clinical	hiv
	surgical	coronary	min	diagnosis	strain
	treatment	ventricular	effect	lesion	infected
	year	myocardial	exercise	brain	patient
	postoperative	cardiac	arterial	syndrome	positive
	operation	left	heart	imaging	viral
k	Topic 16	Topic 17	Topic 18	Topic 19	Topic 20
$n_{k,\bullet}$	22 095 623	21 838 239	21 363 408	20 887 061	20 828 980
PubMed	ca	structure	concentration	pregnancy	protein
	effect	binding	degrees	level	binding
	receptor	protein	samples	women	human
	channel	reaction	liquid	hormone	antibodies
	cell	acid	solution	day	acid
	calcium	interaction	assay	fetal	alpha
	concentration	compound	detection	infant	antibody
	na	site	system	concentration	gel
k	Topic 21	Topic 22	Topic 23	Topic 24	Topic 25
$n_{k,\bullet}$	20 106 260	19 788 488	18 675 096	17 163 327	16 440 018
PubMed	rat	bone	patient	gene	activity
	cell	patient	renal	mutation	acid
	effect	joint	liver	genetic	enzyme
	liver	muscle	transplantati	chromosome	liver
	mice	fractures	blood	analysis	concentration
	dose	hip	disease	genes	rat
	drug	year	acute	dna	enzymes
	mg	implant	chronic	polymorphism	synthesis

k	Topic 26	Topic 27	Topic 28	Topic 29	Topic 30
$n_{k,\bullet}$	16 136 164	14 201 063	13 706 016	13 191 158	13 105 245
PubMed	effect platelet induced oxide rat cell endothelial activity	diet weight intake food body effect acid vitamin	patient disease gastric asthma test pylori arthritis chronic	strain plant growth acid bacteria activity cell species	protein membrane cell domain binding receptor lipid membranes
k	Topic 31	Topic 32	Topic 33	Topic 34	Topic 35
$n_{k,\bullet}$	12 705 261	12 624 252	10 422 885	9 850 167	7 027 660
PubMed	insulin glucose diabetes cholesterol level diabetic plasma lipoprotein	species population infection animal egg host parasite malaria	exposure concentration iron level water effect exposed lead	skin patient eyes eye retinal laser visual corneal	level patient ml control serum plasma factor concentration
k	Topic 36	Topic 37	Topic 38	Topic 39	Topic 40
$n_{k,\bullet}$	6 130 945	644 182	2 264	1 325	104
PubMed	dental oral teeth tooth periodontal treatment salivary gland	sleep caffeine tea effect theophylline night coffee green	ppr csc stretch pthrp response br gei pth	pac foal cpr pacap edm speck branchial lth	feather tieg sorghum coii phycocyanin vanx midrib ifi
k	Topic 41				
$n_{k,\bullet}$	104				
PubMed	steer mca persistency buckwheat dnak eset branding akr				