

Generalists and Specialists

Using Community Embeddings to Quantify
Activity Diversity in Online Platforms

Isaac Waller

walleris@cs.toronto.edu

Ashton Anderson

ashton@cs.toronto.edu

University of Toronto

The Web Conference 2019



Computer Science
UNIVERSITY OF TORONTO

Generalists and specialists

full-stack developer	vs.	React developer
family doctor	vs.	neurosurgeon

Generalists and specialists

full-stack developer	vs.	React developer
family doctor	vs.	neurosurgeon
generalist	vs.	specialist

Generalists and specialists



vulture
generalist



koala
specialist

Koala photo by DAVID ILIFF. License: CC-BY-SA 3.0. Vulture photo by Charles Sharp. License: CC-BY-SA 4.0

Reddit



Games

MakeupAddiction

medicalschoo

soccer

math

programming

Cartalk

chromeos

Construction

funny

television

Aquariums

Which is the specialist?

User 1:

$$C = \{\text{China, nba, Buddhism, startrek}\}$$

User 2:

$$C = \{\text{Fitness, powerlifting, bodybuilding, weightroom}\}$$

Which is the specialist?

User 1:

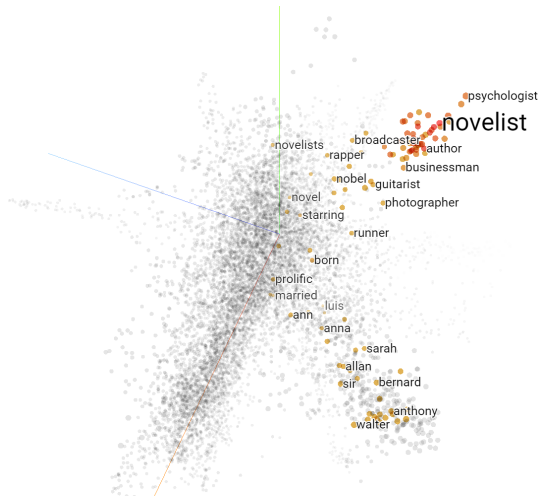
$$C = \{\text{China, nba, Buddhism, startrek}\}$$

User 2:

$$C = \{\text{Fitness, powerlifting, bodybuilding, weightroom}\}$$

$$GS(C) = ?$$

Word2vec¹



[1] Mikolov et al. (2013) Distributed Representations of Words and Phrases and their Compositionality

Word2vec for communities^{2,3}

Input: a (community, user) pair for each comment made in a community

(Games, user1) (Fitness, user3) (medicalschoool, user2)
(China, user4) (Science, user2) (weightlifting, user3)

[2] Kumar et al. (2018) Community Interaction and Conflict on the Web

[3] Martin (2017) community2vec: Vector representations of online communities encode semantic relationships

Word2vec for communities^{2,3}

Input: a (community, user) pair for each comment made in a community

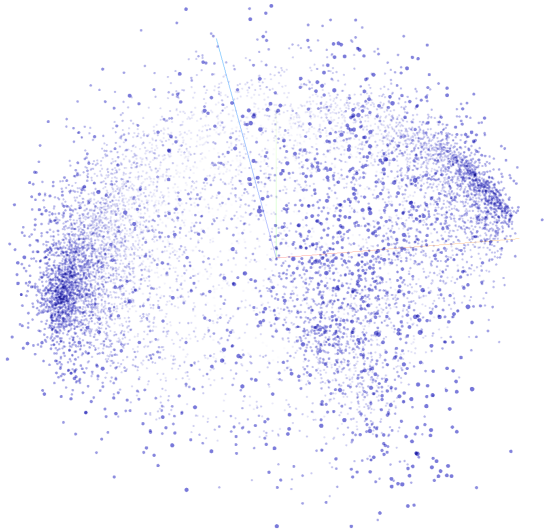
(Games, user1) (Fitness, user3) (medicalschoo, user2)
(China, user4) (Science, user2) (weightlifting, user3)

Output: a vector for each community in the input, where communities with high user overlap are closer to each other

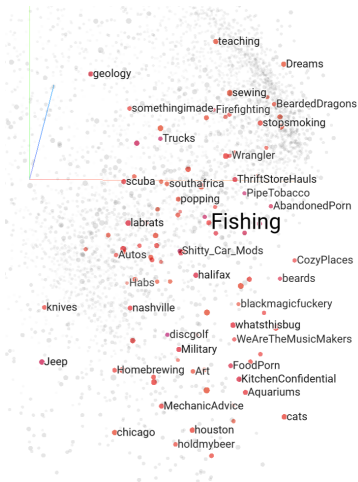
[2] Kumar et al. (2018) Community Interaction and Conflict on the Web

[3] Martin (2017) community2vec: Vector representations of online communities encode semantic relationships

A first embedding



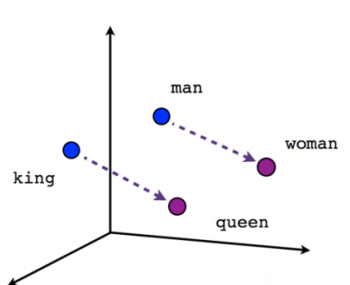
A first embedding



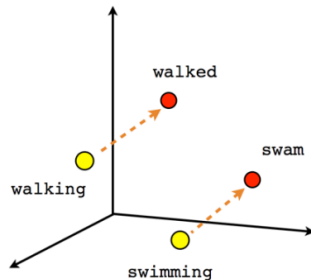
Nearest points in the original space:

discgolf	0.168
Shitty_Car_Mods	0.191
PipeTobacco	0.197
Trucks	0.205
surfing	0.206
flyfishing	0.217
DippingTobacco	0.221
Jeep	0.222
itookapicture	0.236
AbandonedPorn	0.237
beards	0.254
FoodPorn	0.259
halifax	0.261
ThriftStoreHauls	0.265
KitchenConfidential	0.265
labrats	0.270

Word analogies

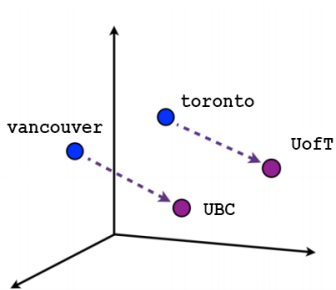


Male to female

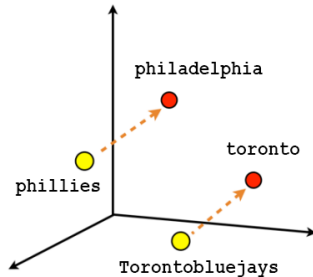


Verb tense

Community analogies



University to city

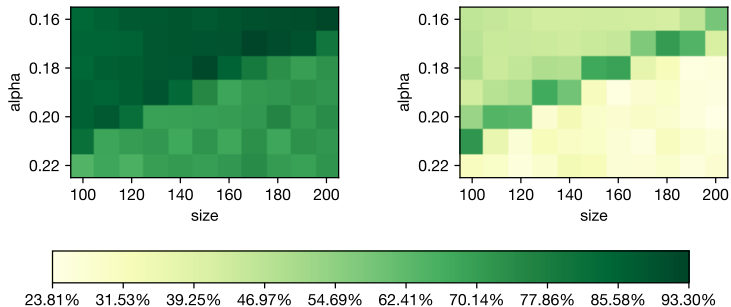


Sports team to sport / city

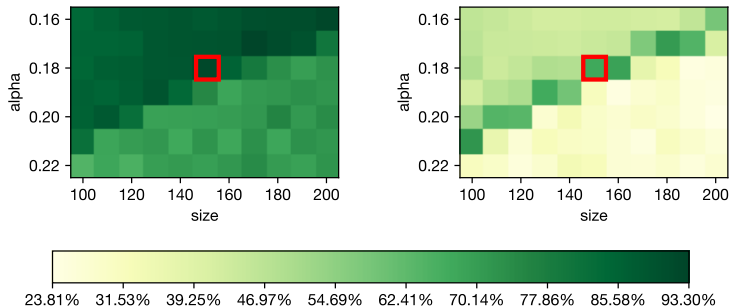
4,392 analogies total

<i>brocku</i>	→	<i>stcatharinesON</i>	as	<i>uakron</i>	→	<i>akron</i>
<i>angelsbaseball</i>	→	<i>baseball</i>	as	<i>LAClippers</i>	→	<i>nba</i>
<i>nus</i>	→	<i>singapore</i>	as	<i>UMT</i>	→	<i>missoula</i>
<i>Colts</i>	→	<i>indianapolis</i>	as	<i>oaklandraiders</i>	→	<i>oakland</i>
<i>PolkStateCollege</i>	→	<i>WinterHaven</i>	as	<i>csun</i>	→	<i>LosAngeles</i>
<i>Coyotes</i>	→	<i>phoenix</i>	as	<i>AnaheimDucks</i>	→	<i>LosAngeles</i>
<i>FLC</i>	→	<i>folsom</i>	as	<i>OxfordBrookes</i>	→	<i>oxford</i>
<i>phillies</i>	→	<i>philadelphia</i>	as	<i>Torontobluejays</i>	→	<i>toronto</i>

Hyperparameter search

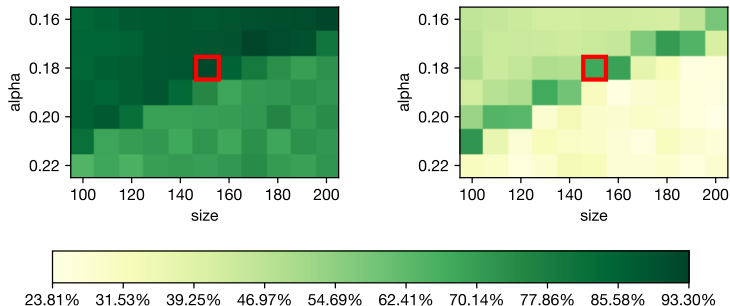


Hyperparameter search



72% perfect, 93% top 5

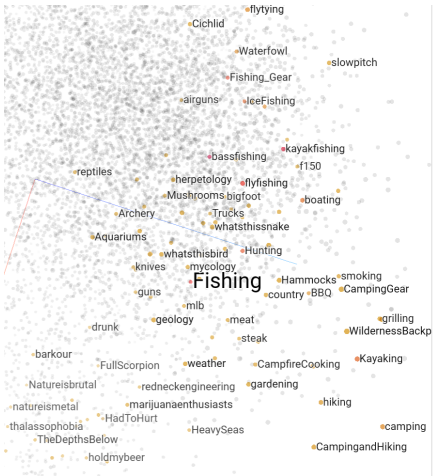
Hyperparameter search



72% perfect, 93% top 5

cycling + swimming + running = triathalon

Our better embedding



Nearest points in the original space:

bassfishing	0.226
kayakfishing	0.260
flyfishing	0.325
Fishing_Gear	0.357
Hunting	0.403
IceFishing	0.440
Kayaking	0.457
boating	0.478
Spearfishing	0.533
flytying	0.542
camping	0.543
Waterfowl	0.548
BBQ	0.562
Outdoors	0.579
bowhunting	0.581
hiking	0.584

Back to generalists and specialists

User 1:

$$C = \{\text{China, nba, Buddhism, startrek}\}$$

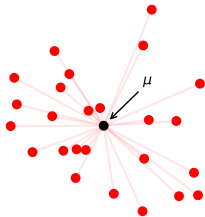
User 2:

$$C = \{\text{Fitness, powerlifting, bodybuilding, weightroom}\}$$

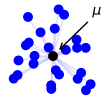
$$GS(C) = ?$$

GS-score

generalist

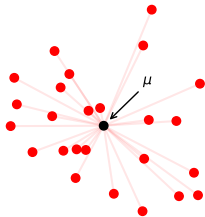


specialist

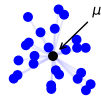


GS-score

generalist



specialist



$$GS(C) = \frac{1}{|C|} \sum_{c \in C} w_c \cos(c, \mu)$$

GS-score

User 1:

$$GS(\{\text{China, nba, Buddhism, startrek}\}) = \frac{0.69}{\mathbf{24^{th} \text{ percentile}}}$$

User 2:

$$GS(\{\text{Fitness, powerlifting, bodybuilding, weightroom}\}) = \frac{0.89}{\mathbf{72^{nd} \text{ percentile}}}$$

$$GS(C) = \frac{1}{|C|} \sum_{c \in C} w_c \cos(c, \mu)$$

Data



All comments in 2017

900M comments, 11.4M distinct users

Top 10,000 subreddits by activity

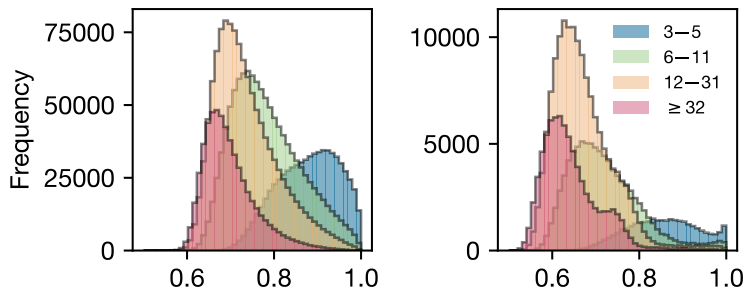


All commits, pull requests, forks,
watches, and stars in 2017

413M actions, 8.3M distinct users

Top 40,000 repos by number of stars

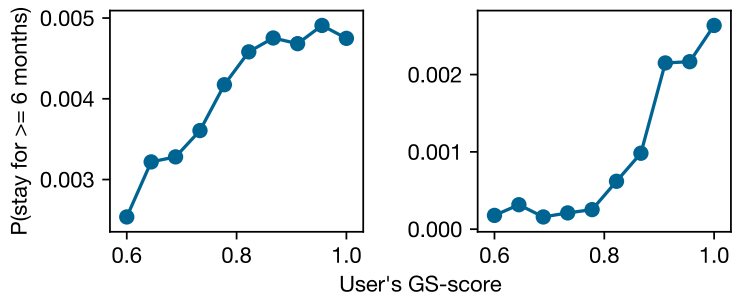
Results



Reddit (left) and GitHub (right)

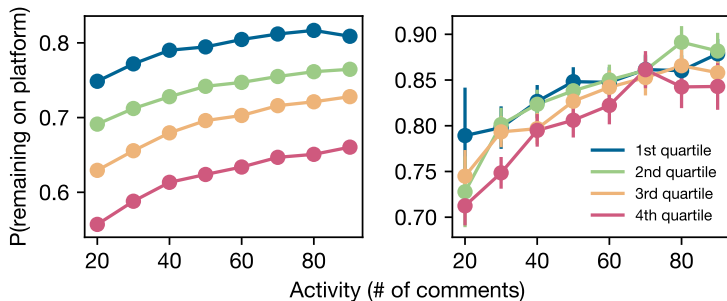
Results

Specialists stay engaged with **communities** longer



Results

Specialists stay engaged with **communities** longer
but generalists stay engaged with the **platform** longer



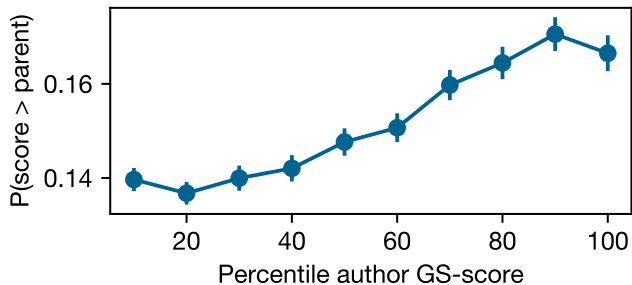
Results

Specialists stay engaged with **communities** longer
but generalists stay engaged with the **platform** longer



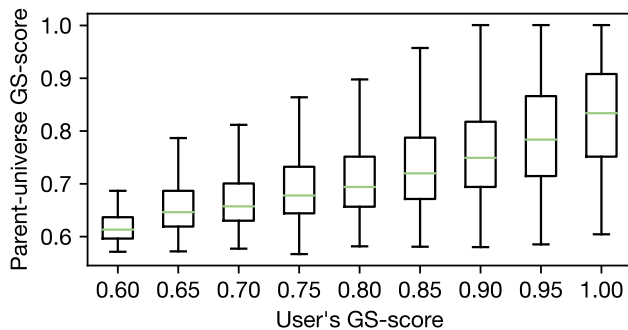
Results

On Reddit, specialists tend to be make more **exceptional comments**



Results

but generalists are exposed to a **more diverse set of users**

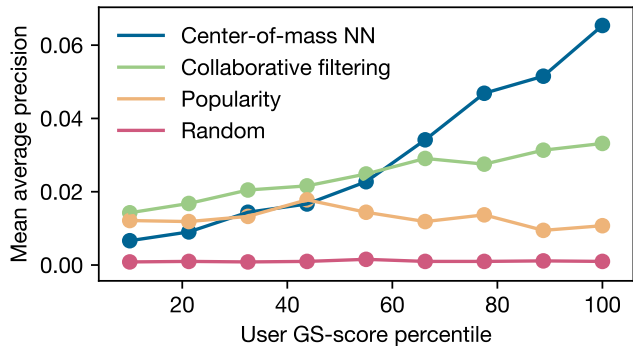


Results

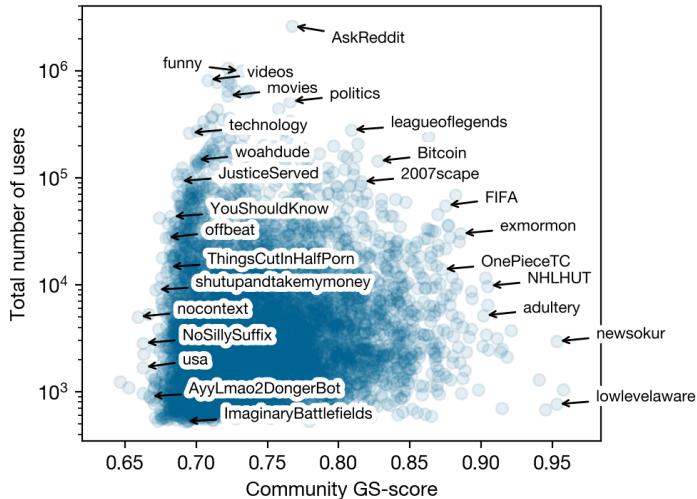
Can GS-score predict new communities a user joins?

Results

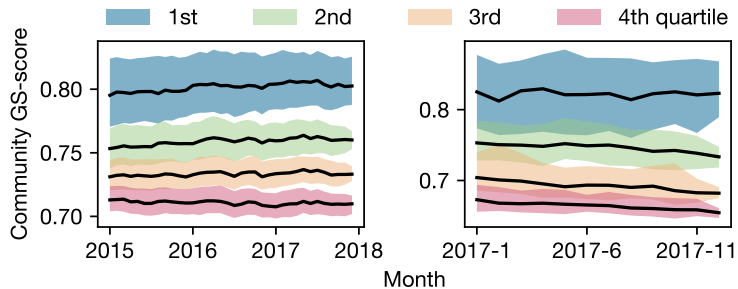
Can GS-score predict new communities a user joins?



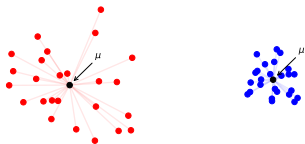
Community GS-scores



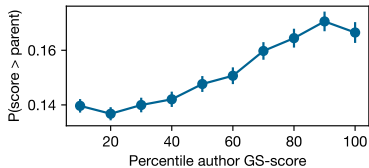
Community GS-scores



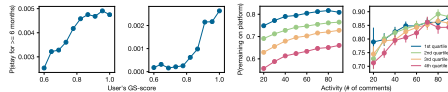
In summary



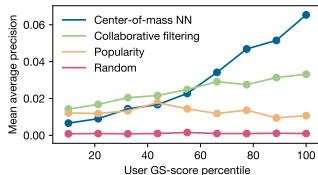
Users on Reddit and GitHub range from generalist to specialist



On Reddit, specialists are more likely to make exceptional comments



Specialists stay engaged with individual communities longer, but generalists stay engaged with the platform longer



Specialists are significantly more predictable than generalists

What does the universe of subreddits look like?

There's a community for almost anything on Reddit. What do these constellations look like? And where are you located within them?

We mapped Reddit to answer this question and find out how users exist in this space. Here's what we found...



A project of the Computational Social Science lab at the University of Toronto ([@isaacwaller](#) and [@ashton1anderson](#)).

▼ Scroll down to continue

Thank you!
tiny.cc/gsscore

