

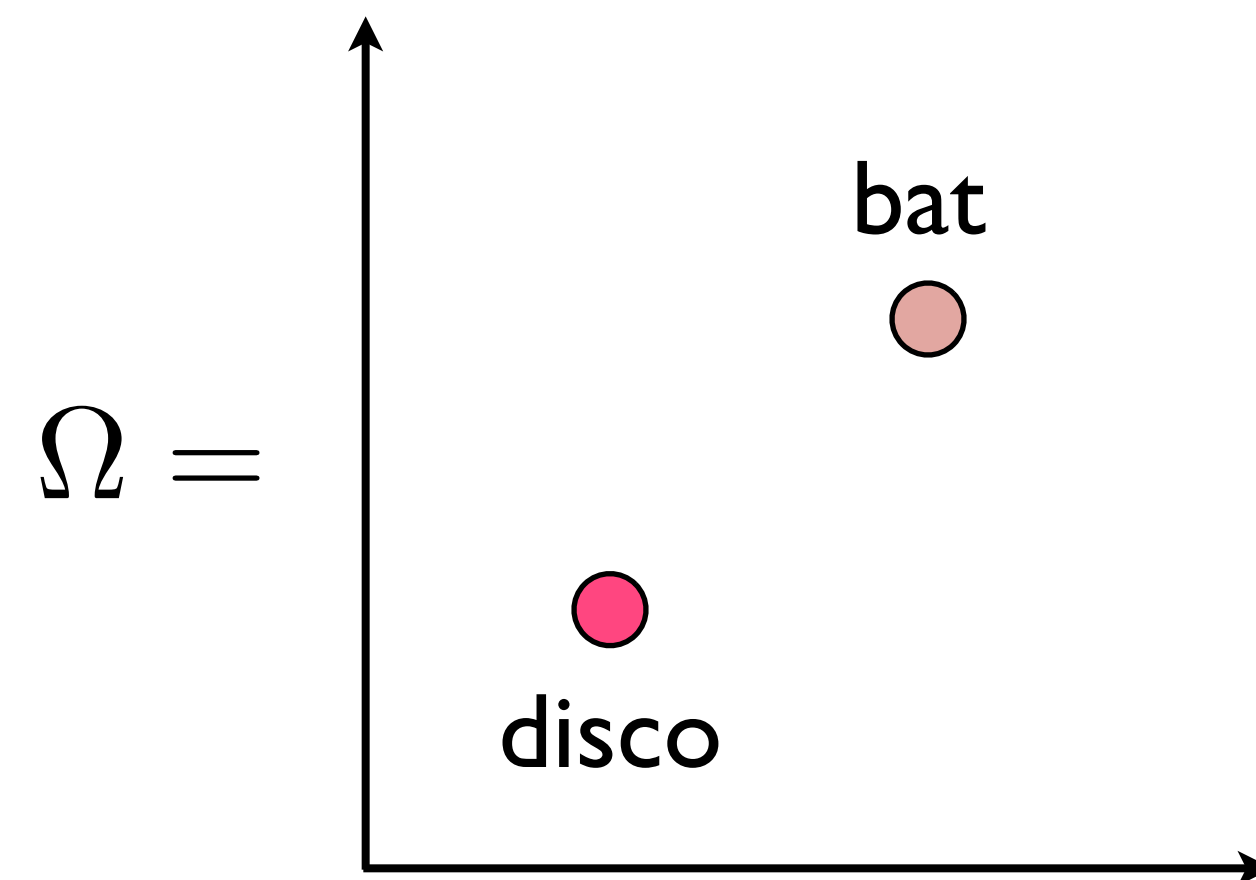
Cross-Cutting Models of Lexical Semantics

Joseph Reisinger and Raymond Mooney

Distributional Lexical Semantics

- Represent “meaning” as a point/vector in a high-dimensional space
- Word relatedness correlates with some distance metric

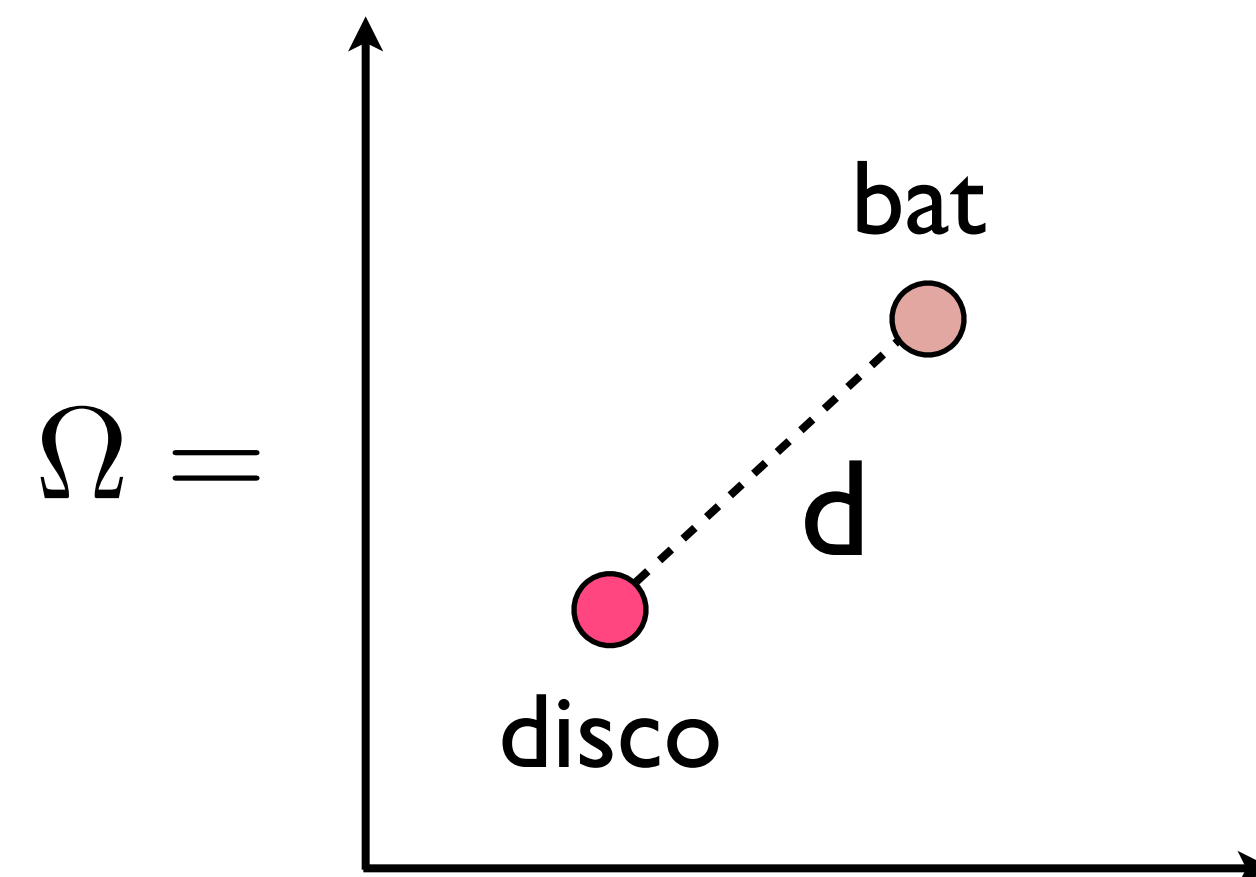
Almuhareb and Poesio (2004), Baroni and Lenci (2009), Bullinaria and Levy (2007), Erk (2007), Griffiths et al. (2007), Landauer and Dumais (1997), Moldovan (2006), Padó and Lapata (2007), Pantel and Pennacchiotti (2006), Sahlgren (2006), Turney and Pantel (2010)



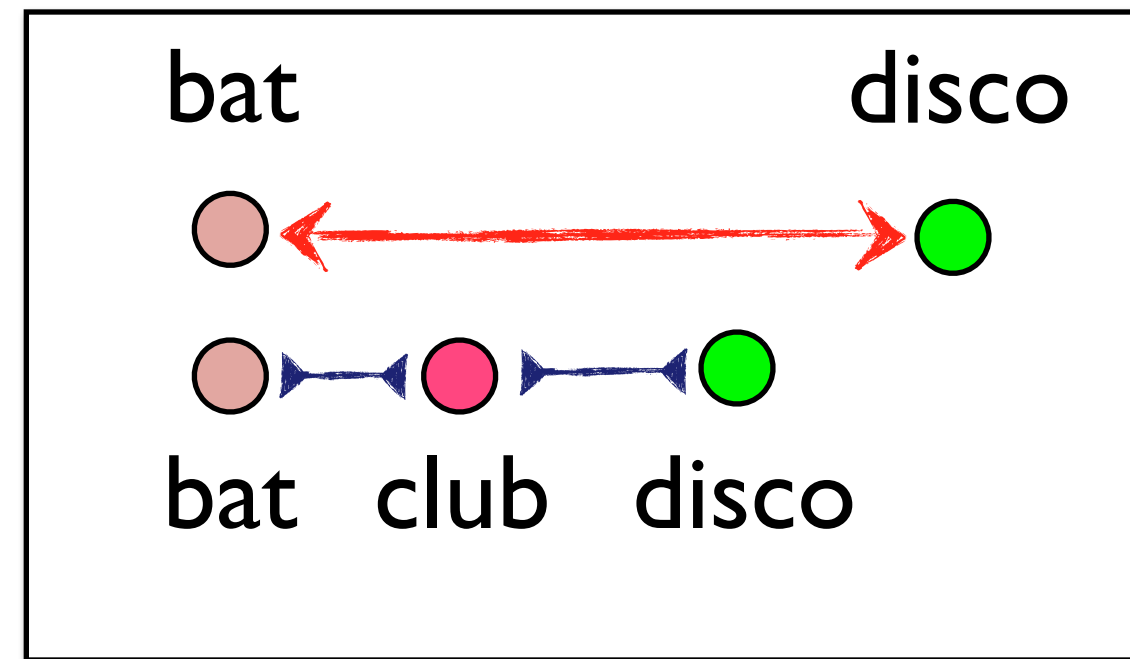
Distributional Lexical Semantics

- Represent “meaning” as a point/vector in a high-dimensional space
- Word relatedness correlates with some distance metric

Almuhareb and Poesio (2004), Baroni and Lenci (2009), Bullinaria and Levy (2007), Erk (2007), Griffiths et al. (2007), Landauer and Dumais (1997), Moldovan (2006), Padó and Lapata (2007), Pantel and Pennacchiotti (2006), Sahlgren (2006), Turney and Pantel (2010)



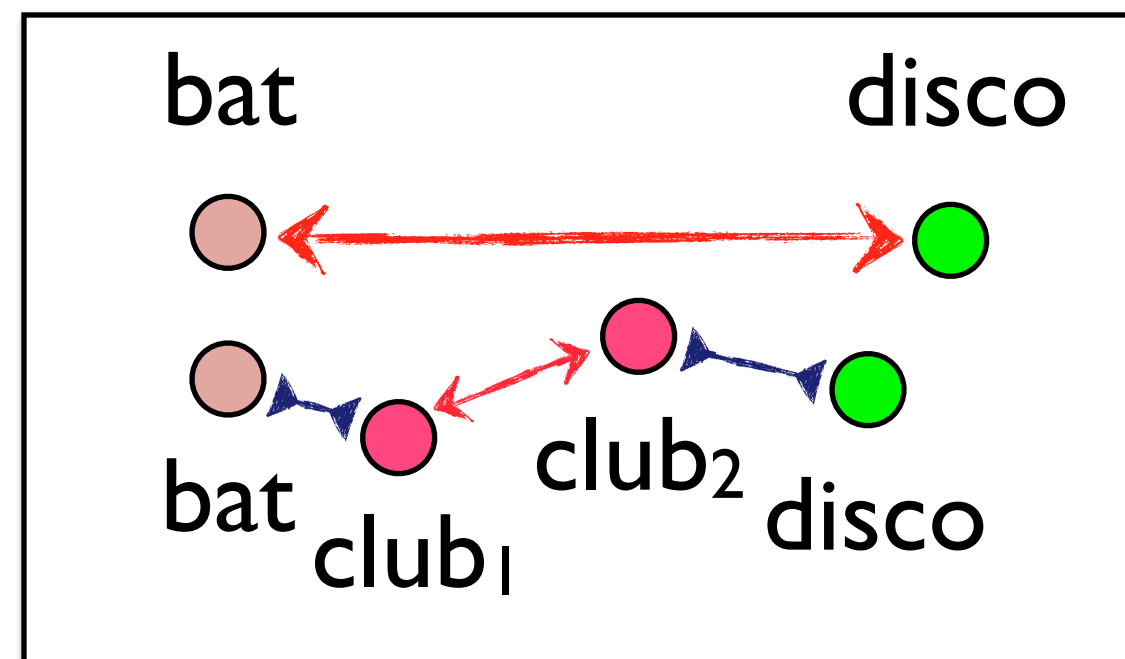
Distributional Lexical Semantics



“meaning violates the triangle inequality”

Tversky and Gati (1982), Griffiths et al. (2007)

Distributional Lexical Semantics

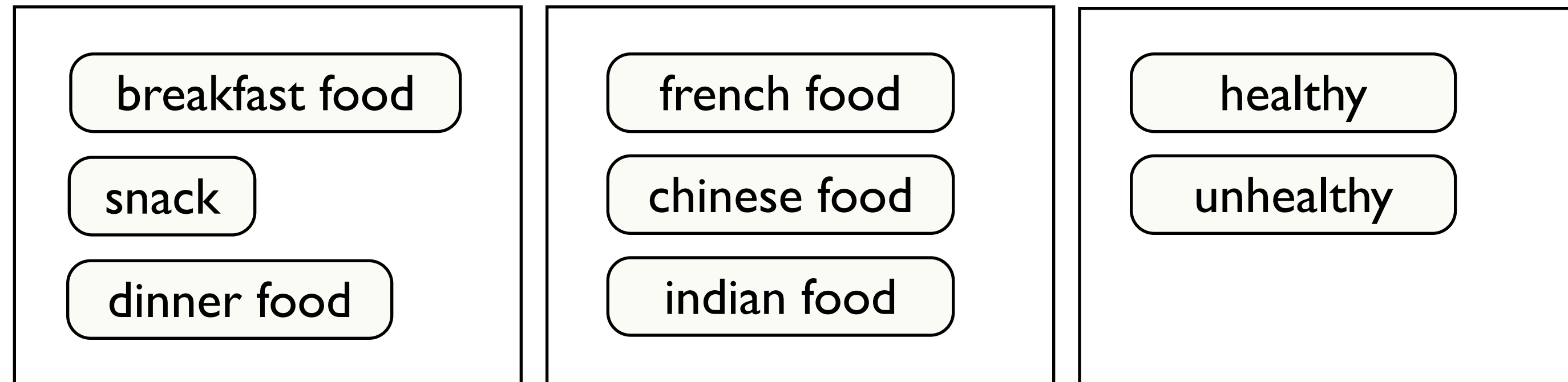


“meaning violates the triangle inequality”

Tversky and Gati (1982), Griffiths et al. (2007)

- Address metric violations by learning word sense clusters / making use of local context
- Can we build a model that captures this directly?

Cross-cutting Concept Organization



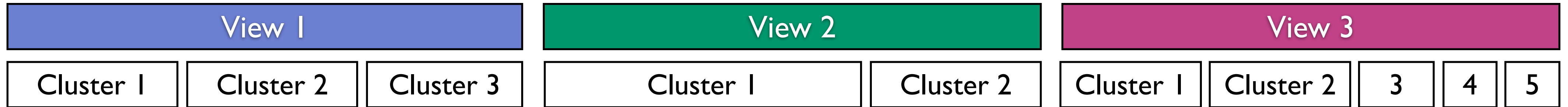
- Human concept organization exhibits **cross cutting** structure
Rosch, et al. (1976); Ross & Murphy (1999); Medin, et al. (2005); Shaftoe, et al. (2011)
- Each categorization system controls what kinds of generalizations (e.g. inferences) are valid.
- Do word usages exhibit similar cross-cutting?
- *Xue, Chen and Palmer (2006)*: sense disambiguation requires vastly different features for different polysemous verbs in Chinese.

Multi View Multinomial Clustering

- There are many valid word clusterings, each capturing different aspects of syntax or topicality
- We introduce a model to explicitly capture multiple organizational systems
- **Cross-cutting** categorization / latent subspaces with separate, coherent clusterings
- Implement using LDA and DPMM primitives / Gibbs sampling

<p>and is ____ and are ____ we are ____ which was ____ he is ____ who are ____</p>		
<p>unwilling willing reluctant refusing glad</p>	<p>exceedingly sincerely logically justly appropriately</p>	<p>about because</p>
<p>brand new ____ results for ____ selection of ____ the latest ____ ____ for sale to buy ____</p>		
<p>samsung panasonic toshiba sony epson</p>	<p>toyota nissan mercedes volvo audi</p>	<p>dunlop yokohama toyo uniroyal michelin</p>

Multi View Multinomial Clustering Model



Data

Austin

History of Austin, Texas, University of Texas Medical Branch, 1993 Pacific hurricane season, Rutherford B. Hayes, List of pipeline accidents, List of Austin City Limits performers, Texas in the American Civil War, 6th Cavalry Regiment (United States)
 ___ texas homes, ___ law school, the citizens of ___, the ___ business directory, ___ police department, university in ___,
 ___ vacation rentals, the ___ parks and, by the ___ business journal, coming to ___, the ___ area, deals on ___ hotels

Betrayed

Survivor: The Amazon, Personal life of Marcus Tullius Cicero, Numb3rs, Huns, Rurouni Kenshin, Liberation of Paris, The Knightly Tale of Gologras and Gawain, Territories in The Pendragon Adventure, A Storm of Swords, Connor MacLeod, Paul Atreides
 her manner ___, being ___ by their, ___ and murdered, ___ his weakness, she ___ him, ___ the secret, ___ by her husband, a voice that ___, who felt ___, ___ to the police, ___ their country, suspected of having ___, ___ the confidence, even when ___

Cat

South China Tiger, Hybrid (biology), List of mammals of Cameroon, Cantonese cuisine, Pound Puppies, Wonder Pets, The Wizard of Oz (1902 stage play), Mee-Ow, Animal rights, Rickrolling, Mera (comics), Taboo food and drink, Tuna, Garfield: The Movie
 ate the ___, have a ___ and a, the ___ and the mouse, the ___ who killed, ___ toys by, ___ in the city, ___ was diagnosed, crazy ___ lady, ___ of the month, protect your ___ from, new ___ food, and bought a ___, ___ or other animal, a sick ___,

View 1

Cluster 1

Cluster 2

Cluster 3

View 2

Cluster 1

Cluster 2

View 3

Cluster 1

Cluster 2

3

4

5

View 1

Cluster 1

Cluster 2

Cluster 3

View 2

Cluster 1

Cluster 2

View 3

Cluster 1

Cluster 2

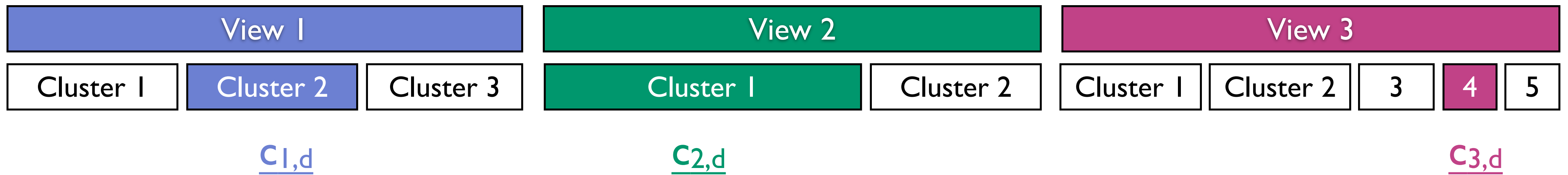
3

4

5

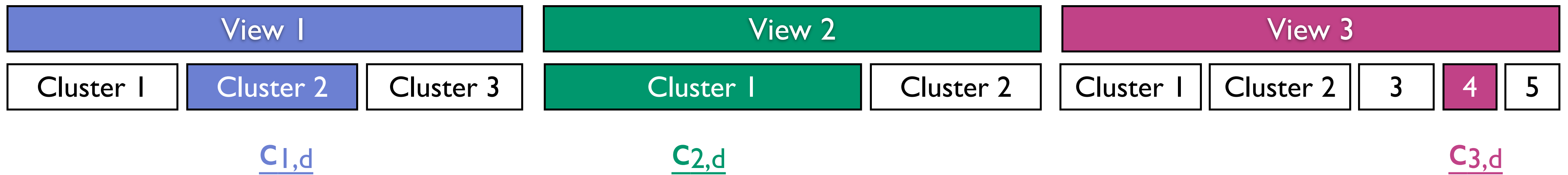
Cat

South China Tiger, Hybrid (biology), List of mammals of Cameroon, Cantonese cuisine, Pound Puppies, Wonder Pets, The Wizard of Oz (1902 stage play), Mee-Ow, Animal rights, Rickrolling, Mera (comics), Taboo food and drink, Tuna, Garfield: The Movie ate the ____, have a ____ and a, the ____ and the mouse, the ____ who killed, ____ toys by, ____ in the city, ____ was diagnosed, crazy ____ lady, ____ of the month, protect your ____ from, new ____ food, and bought a ____, ____ or other animal, a sick ____,



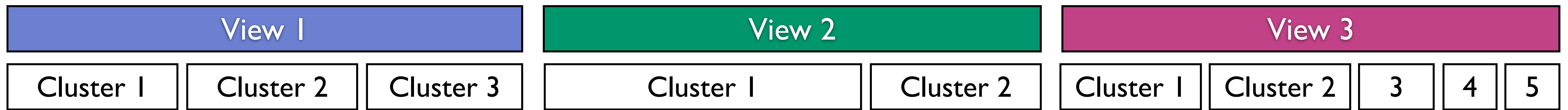
Cat South China Tiger, Hybrid (biology), List of mammals of Cameroon, Cantonese cuisine, Pound Puppies, Wonder Pets, The Wizard of Oz (1902 stage play), Mee-Ow, Animal rights, Rickrolling, Mera (comics), Taboo food and drink, Tuna, Garfield: The Movie ate the ____, have a ____ and a, the ____ and the mouse, the ____ who killed, ____ toys by, ____ in the city, ____ was diagnosed, crazy ____ lady, ____ of the month, protect your ____ from, new ____ food, and bought a ____, ____ or other animal, a sick ____.

- Select a cluster assignment $c_{v,d}$ for d in each view v (DPMM) **i.e. words are assigned to clusters within each view**



Cat South China Tiger, Hybrid (biology), List of mammals of Cameroon, Cantonese cuisine, Pound Puppies, Wonder Pets, The Wizard of Oz (1902 stage play), Mee-Ow, Animal rights, Rickrolling, Mera (comics), Taboo food and drink, Tuna, Garfield: The Movie ate the ____, have a ____ and a, the ____ and the mouse, the ____ who killed, ____ toys by, ____ in the city, ____ was diagnosed, crazy ____ lady, ____ of the month, protect your ____ from, new ____ food, and bought a ____, ____ or other animal, a sick ____.

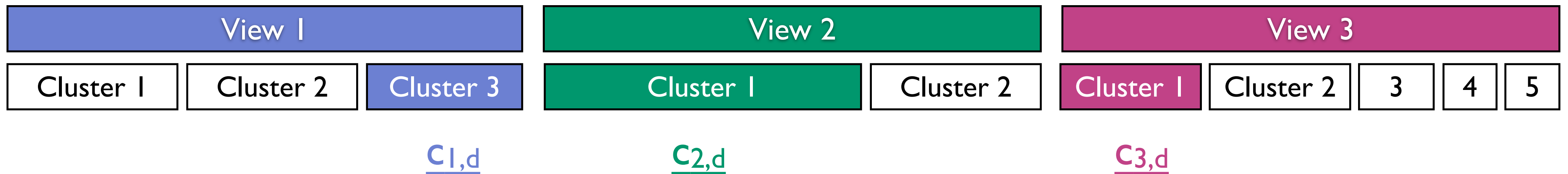
- Select a cluster assignment $c_{v,d}$ for d in each view v (DPMM) *i.e. words are assigned to clusters within each view*
- Select a view v_f for each observed feature, and generate it from $c_{v_f,d}$ (LDA) *i.e. features distributed between views*



Cat South China Tiger, Hybrid (biology), List of mammals of Cameroon, Cantonese cuisine, Pound Puppies, Wonder Pets, The Wizard of Oz (1902 stage play), Mee-Ow, Animal rights, Rickrolling, Mera (comics), Taboo food and drink, Tuna, Garfield: The Movie
 ate the ____, have a ____ and a, the ____ and the mouse, the ____ who killed, ____ toys by, ____ in the city, ____ was diagnosed, crazy ____ lady, ____ of the month, protect your ____ from, new ____ food, and bought a ____, ____ or other animal, a sick ____,

Betrayed Survivor: The Amazon, Personal life of Marcus Tullius Cicero, Numb3rs, Huns, Rurouni Kenshin, Liberation of Paris, The Knightly Tale of Gologras and Gawain, Territories in The Pendragon Adventure, A Storm of Swords, Connor MacLeod, Paul Atreides
 her manner ____, being ____ by their, ____ and murdered, ____ his weakness, she ____ him, ____ the secret, ____ by her husband, a voice that ____, who felt ____, ____ to the police, ____ their country, suspected of having ____, ____ the confidence, even when ____

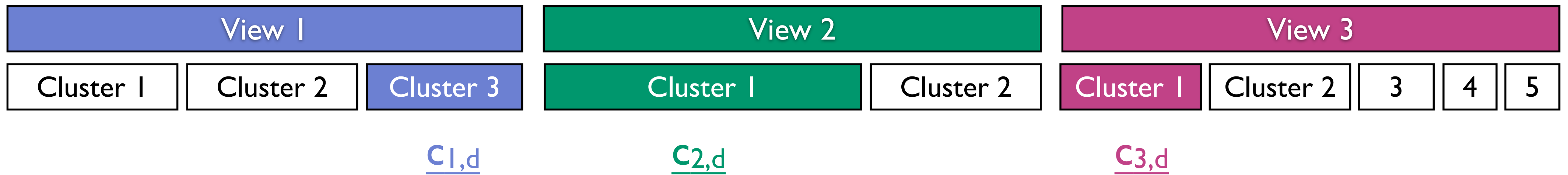
- Select a cluster assignment $c_{v,d}$ for d in each view v (DPMM) **i.e. words are assigned to clusters within each view**
- Select a view v_f for each observed feature, and generate it from $c_{v_f,d}$ (LDA) **i.e. features distributed between views**



Cat South China Tiger, Hybrid (biology), List of mammals of Cameroon, Cantonese cuisine, Pound Puppies, Wonder Pets, The Wizard of Oz (1902 stage play), Mee-Ow, Animal rights, Rickrolling, Mera (comics), Taboo food and drink, Tuna, Garfield: The Movie
 ate the ____, have a ____ and a, the ____ and the mouse, the ____ who killed, ____ toys by, ____ in the city, ____ was diagnosed, crazy ____ lady, ____ of the month, protect your ____ from, new ____ food, and bought a ____, ____ or other animal, a sick ____,

Betrayed Survivor: The Amazon, Personal life of Marcus Tullius Cicero, Numb3rs, Huns, Rurouni Kenshin, Liberation of Paris, The Knightly Tale of Gologras and Gawain, Territories in The Pendragon Adventure, A Storm of Swords, Connor MacLeod, Paul Atreides
 her manner ____, being ____ by their, ____ and murdered, ____ his weakness, she ____ him, ____ the secret, ____ by her husband, a voice that ____, who felt ____, ____ to the police, ____ their country, suspected of having ____, ____ the confidence, even when ____

- Select a cluster assignment $c_{v,d}$ for d in each view v (DPMM) **i.e. words are assigned to clusters within each view**
- Select a view v_f for each observed feature, and generate it from $c_{v_f,d}$ (LDA) **i.e. features distributed between views**



Cat

South China Tiger, Hybrid (biology), List of mammals of Cameroon, Cantonese cuisine, Pound Puppies, Wonder Pets, The Wizard of Oz (1902 stage play), Mee-Ow, Animal rights, Rickrolling, Mera (comics), Taboo food and drink, Tuna, Garfield: The Movie ate the ____, have a ____ and a, the ____ and the mouse, the ____ who killed, ____ toys by, ____ in the city, ____ was diagnosed, crazy ____ lady, ____ of the month, protect your ____ from, new ____ food, and bought a ____, ____ or other animal, a sick ____,

Betrayed

Survivor: The Amazon, Personal life of Marcus Tullius Cicero, Numb3rs, Huns, Rurouni Kenshin, Liberation of Paris, The Knightly Tale of Gologras and Gawain, Territories in The Pendragon Adventure, A Storm of Swords, Connor MacLeod, Paul Atreides her manner ____, being ____ by their, ____ and murdered, ____ his weakness, she ____ him, ____ the secret, ____ by her husband, a voice that ____, who felt ____, ____ to the police, ____ their country, suspected of having ____, ____ the confidence, even when ____

- Select a cluster assignment $c_{v,d}$ for d in each view v (DPMM) **i.e. words are assigned to clusters within each view**
- Select a view v_f for each observed feature, and generate it from $c_{v_f,d}$ (LDA) **i.e. features distributed between views**

View 1

View 2

View 3

Cluster 1

Cluster 2

Cluster 1

Cluster 1

Cluster 2

to an
side of the
first of
of human
the little
of from the
and an
written by
way of
real estate in
of may
hotels hotels
estate in
city of
welcome to
was the of
town of
to a
the city of
the does not
private message to
presence of
posted by at
name of
message to
located in
like and
in an
in the
hotels in
going to
from the to
dsl dsl
degree of
create a
by to
by on
born in
an and
was born
said that
high school
do not
you are
who is
which was
which the
were in
we are
was to
to be and
the more
so many
she was
of the were
of a and
near the
is also
i was
his of
could be
been and
be or
as as
and was
and is
and are
and his
also the
a more
some of
who are
were not
the very
the american
the of that
the must be
the family
that was
that are
posts by
of being
of have
might be
many and
is an
in these
he is
but the of
be to
are to
and their
along the
a kind of
who had
open this result in
home page

arbitrary
austin
batimore
characteristic
comparative
dallas
evolutionary
franklin
fundamental
inadequate
integral
jackson
kent
likeilhood
liverpool
mystical
newcastle
pittsburgh
poetic
proportional
psychological
radical
richmond
singular

betrayed
conquered
disappointed
divorced
embarked
frustrated
guaraded
hated
knocked
murdered
praised
stationed
stole
summoned
wounded
secretly

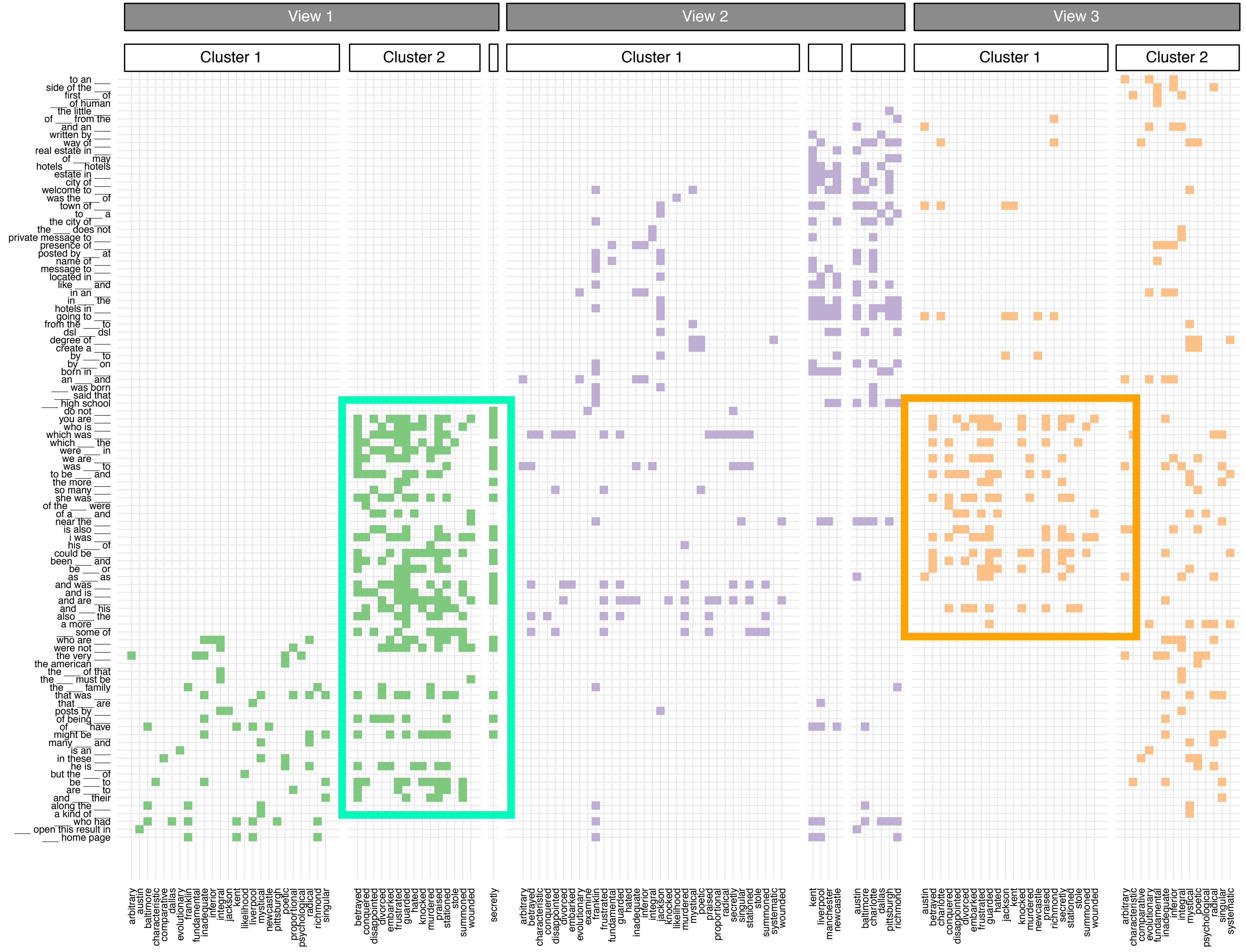
arbitrary
betrayed
characteristic
conquered
disappointed
divorced
embarked
evolutionary
examine
franklin
frustrated
fundamental
guaraded
hated
inadequate
interior
integral
jackson
knocked
likeilhood
murdered
mystical
poetic
praised
proportional
radical
secretly
singular
stationed
stole
summoned
systematic
wounded

kent
liverpool
manchester
newcastle

austin
batimore
charlotte
dallas
pittsburgh
richmond

austin
charlotte
conquered
disappointed
embarked
frustrated
guaraded
hated
jackson
kent
knocked
murdered
newcastle
praised
richmond
secretly
stationed
stole
summoned
wounded

arbitrary
characteristic
comparative
evolutionary
fundamental
inadequate
interior
integral
mystical
poetic
psychological
radical
singular
systematic



to an
side of the
first of
of human
the little
of from the
and an
written by
way of
real estate in
of may
hotels hotels
estate in
city of
welcome to
was the of
town of
to a
the city of
the does not
private message to
presence of
posted by at
name of
message to
located in
like and
in an
in the
hotels in
going to
from the to
dsl dsl
degree of
create a
by to
by on
born in
an and
was born
said that
high school
do not
you are
who is
which was
which the
were in
we are
was to
to be and
the more
so many
she was
of the were
of a and
near the
is also
i was
his of
could be
been and
be or
as as
and was
and is
and are
and his
also the
a more
some of
who are
were not
the very
the american
the of that
the must be
the family
that was
that are
posts by
of being
of have
might be
many and
is an
in these
he is
but the of
be to
are to
and their
along the
a kind of
who had
open this result in
home page

arbitrary
austin
balTIMORE
characteristic
comparative
dallas
evolutionary
franklin
fundamental
inadequate
integral
jackson
kent
likeIhood
liverpool
mystical
newcastle
pittsburgh
poetic
proportional
psychological
radical
richmond
singular

betrayed
conquered
disappointed
divorced
embarked
frustrated
guaranteed
hated
knocked
murdered
praised
stationed
stole
summoned
wounded
secretly

arbitrary
betrayed
characteristic
conquered
disappointed
divorced
embarked
frustrated
evolutionary
examine
franklin
frustrated
fundamental
guaranteed
hated
inadequate
interior
integral
jackson
knocked
likeIhood
murdered
mystical
poetic
proportional
radical
secretly
singular
stationed
stole
summoned
systematic
wounded

kent
liverpool
manchester
newcastle
austin
balTIMORE
charlotte
dallas
pittsburgh
richmond
austin
charlotte
conquered
disappointed
divorced
embarked
frustrated
guaranteed
hated
jackson
kent
knocked
murdered
newcastle
praised
richmond
secretly
stationed
stole
summoned
wounded
arbitrary
characteristic
comparative
evolutionary
fundamental
inadequate
interior
integral
mystical
poetic
psychological
radical
singular
systematic

Data

Austin

History of Austin, Texas, University of Texas Medical Branch, 1993 Pacific hurricane season, Rutherford B. Hayes, List of pipeline accidents, List of Austin City Limits performers, Texas in the American Civil War, 6th Cavalry Regiment (United States)
___ texas homes, ___ law school, the citizens of ___, the ___ business directory, ___ police department, university in ___, ___ vacation rentals, the ___ parks and, by the ___ business journal, coming to ___, the ___ area, deals on ___ hotels

Betrayed

Survivor: The Amazon, Personal life of Marcus Tullius Cicero, Numb3rs, Huns, Rurouni Kenshin, Liberation of Paris, The Knightly Tale of Gologras and Gawain, Territories in The Pendragon Adventure, A Storm of Swords, Connor MacLeod, Paul Atreides
her manner ___, being ___ by their, ___ and murdered, ___ his weakness, she ___ him, ___ the secret, ___ by her husband, a voice that ___, who felt ___, ___ to the police, ___ their country, suspected of having ___, ___ the confidence, even when ___

- **Word set:** Top 43.7k words ranked by frequency in Wikipedia (ex top 1% as stop words)
- **Syntax features:** Contextual patterns from combined Google Web n-gram + Google Books n-gram corpus (3.5M features)
- **Document features:** Wikipedia article occurrence count (120k features)

Intrusion Task

word

humor
ingenuity
delight
advertisers
astonishment

context

document

- “Model-based” lexical semantics: read word similarity directly from the model
- Intruders are drawn from the top terms in other clusters
- More robust than asking for numeric similarity judgements
- Less inter-rater calibration required

Intrusion Task

word

context

document

humor
ingenuity
delight
advertisers
astonishment

- “Model-based” lexical semantics: read word similarity directly from the model
- Intruders are drawn from the top terms in other clusters
- More robust than asking for numeric similarity judgements
- Less inter-rater calibration required

Intrusion Task

word

humor
ingenuity
delight
advertisers
astonishment

context

_____ is characterized
symptoms of _____
cases of _____
in cases of _____
real estate in _____

document

- “Model-based” lexical semantics: read word similarity directly from the model
- Intruders are drawn from the top terms in other clusters
- More robust than asking for numeric similarity judgements
- Less inter-rater calibration required

Intrusion Task

word

humor
ingenuity
delight
advertisers
astonishment

context

_____ is characterized
symptoms of _____
cases of _____
in cases of _____
real estate in _____

document

- “Model-based” lexical semantics: read word similarity directly from the model
- Intruders are drawn from the top terms in other clusters
- More robust than asking for numeric similarity judgements
- Less inter-rater calibration required

Intrusion Task

word

humor
ingenuity
delight
advertisers
astonishment

context

_____ is characterized
symptoms of _____
cases of _____
in cases of _____
real estate in _____

document

Puerto Rican cuisine
Greek cuisine
ThinkPad
Palestinian cuisine
Field ration

- “Model-based” lexical semantics: read word similarity directly from the model
- Intruders are drawn from the top terms in other clusters
- More robust than asking for numeric similarity judgements
- Less inter-rater calibration required

Intrusion Task

word

humor
ingenuity
delight
advertisers
astonishment

context

_____ is characterized
symptoms of _____
cases of _____
in cases of _____
real estate in _____

document

Puerto Rican cuisine
Greek cuisine
ThinkPad
Palestinian cuisine
Field ration

- “Model-based” lexical semantics: read word similarity directly from the model
- Intruders are drawn from the top terms in other clusters
- More robust than asking for numeric similarity judgements
- Less inter-rater calibration required

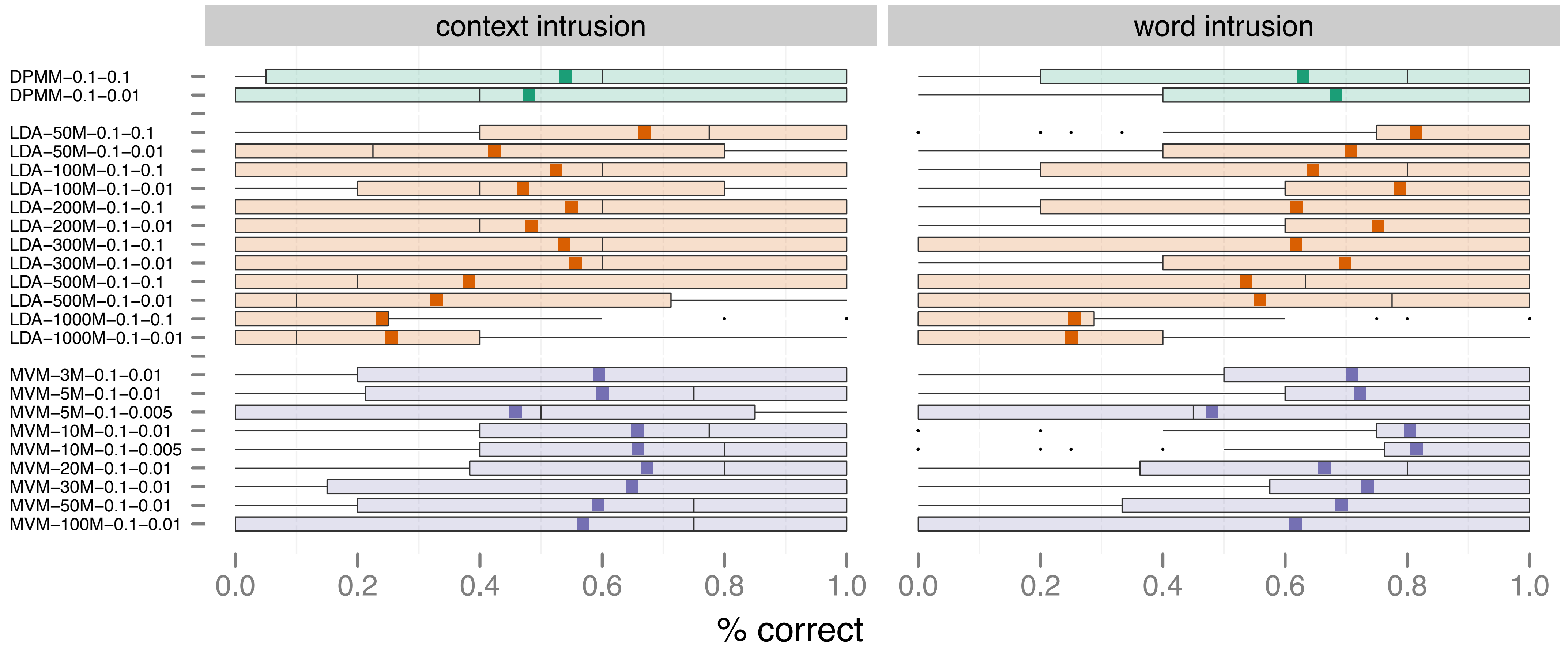
Evaluation

- Amazon Mechanical Turk
- 1256 unique raters (Country=US, >96% approval)
- 5.7k unique intrusion tasks at 5x duplication: ~30k evaluations total
- 2736 rejected
 - Per-user average time for <1.5s / question
 - Low-entropy answers
 - Low agreement

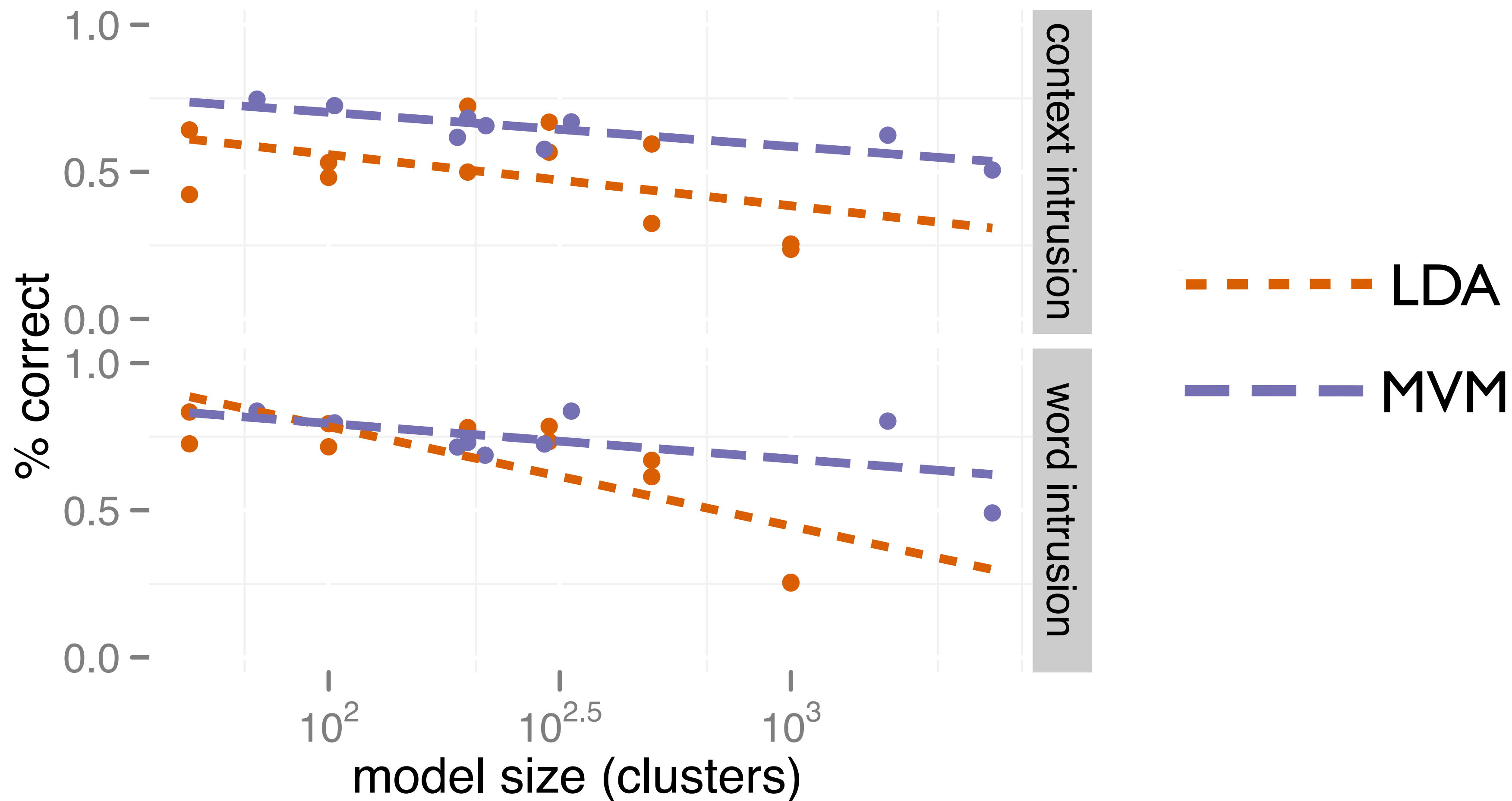
User Comments

- U1 I just tried 30 of the what doesn't belong ones. They took about 30 seconds each due to thinking time so not worth it for me.
 - U2 I don't understand the fill in the blank ones to be honest. I just kinda pick one, since I don't know what's expected lol
 - U3 Your not filling in the blank just ignore the blank and think about how the words they show relate to each other and choose the one that relates least. Some have just words and no blanks.
 - U4 These seem very subjective to mw. i hope there isn't definite correct answers because some of them make me go [emoticon of head-scratching]
 - U5 I looked and have no idea. I guess I'm a word idiot because I don't see the relation between the words in the preview HIT - too scared to try any of these.
 - U6 I didn't dive in but I did more than I should have they were just too easy. Most of them I could tell what did not belong, some were pretty iffy though.
-

Syntax features only (freq>50; “common”)



Syntax features only (freq>50; “common”)

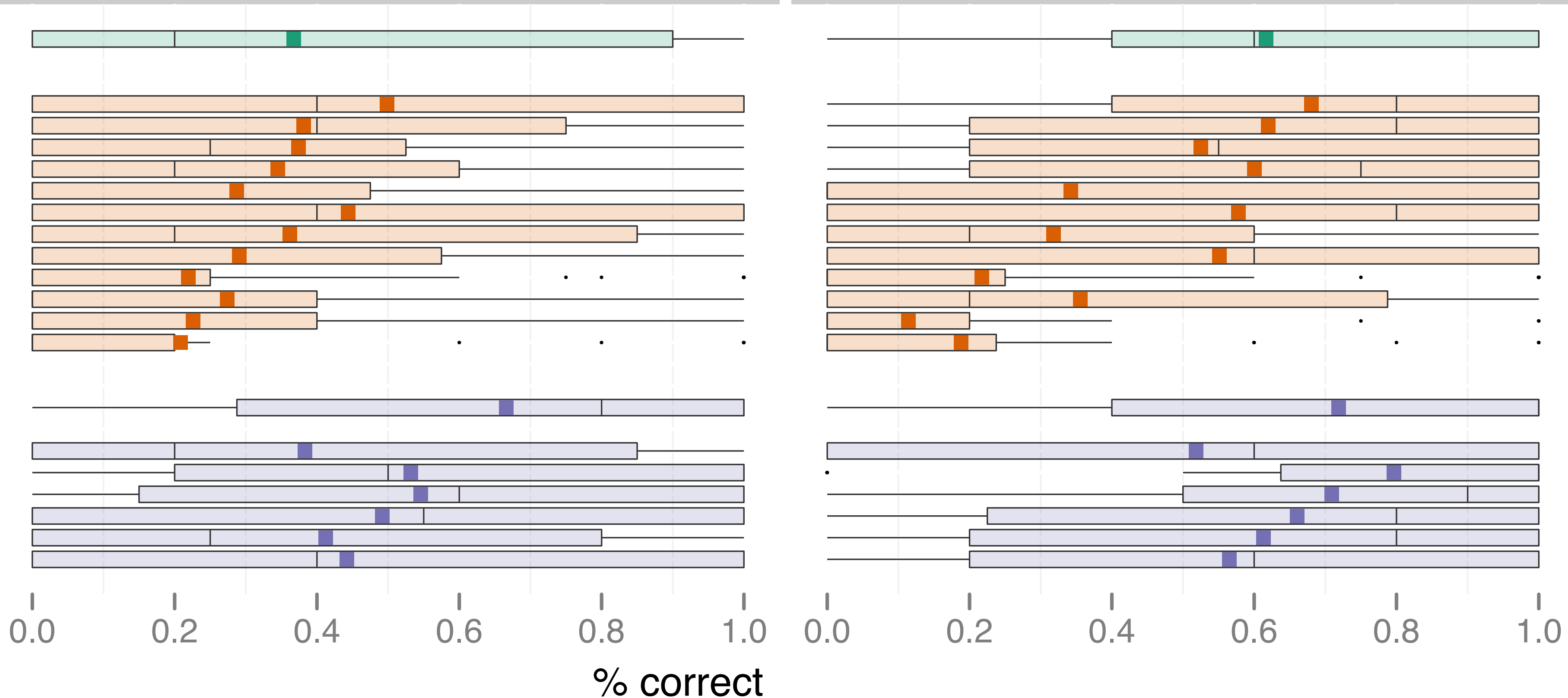


Syntax features only (freq < 50; “rare”)

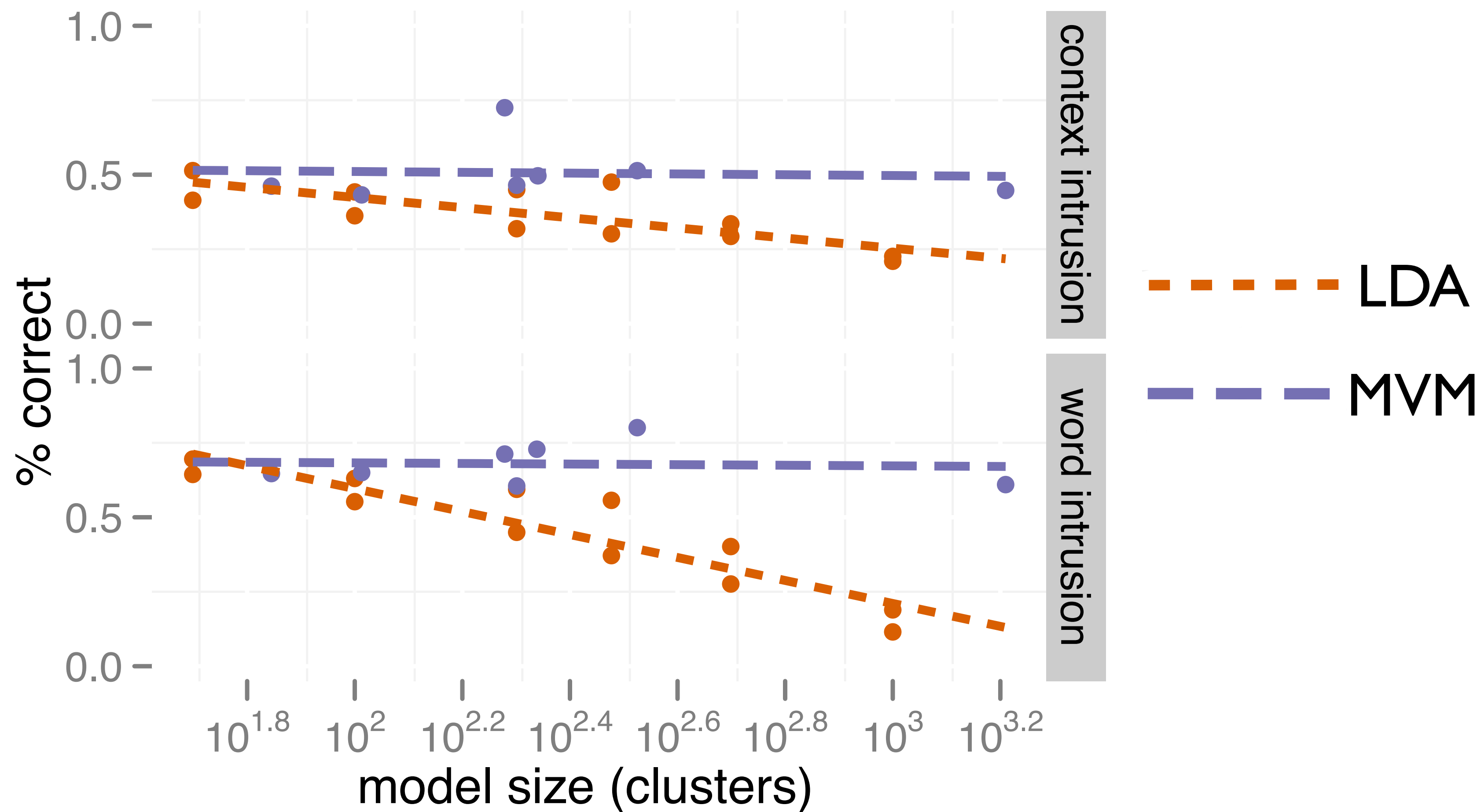
context intrusion

word intrusion

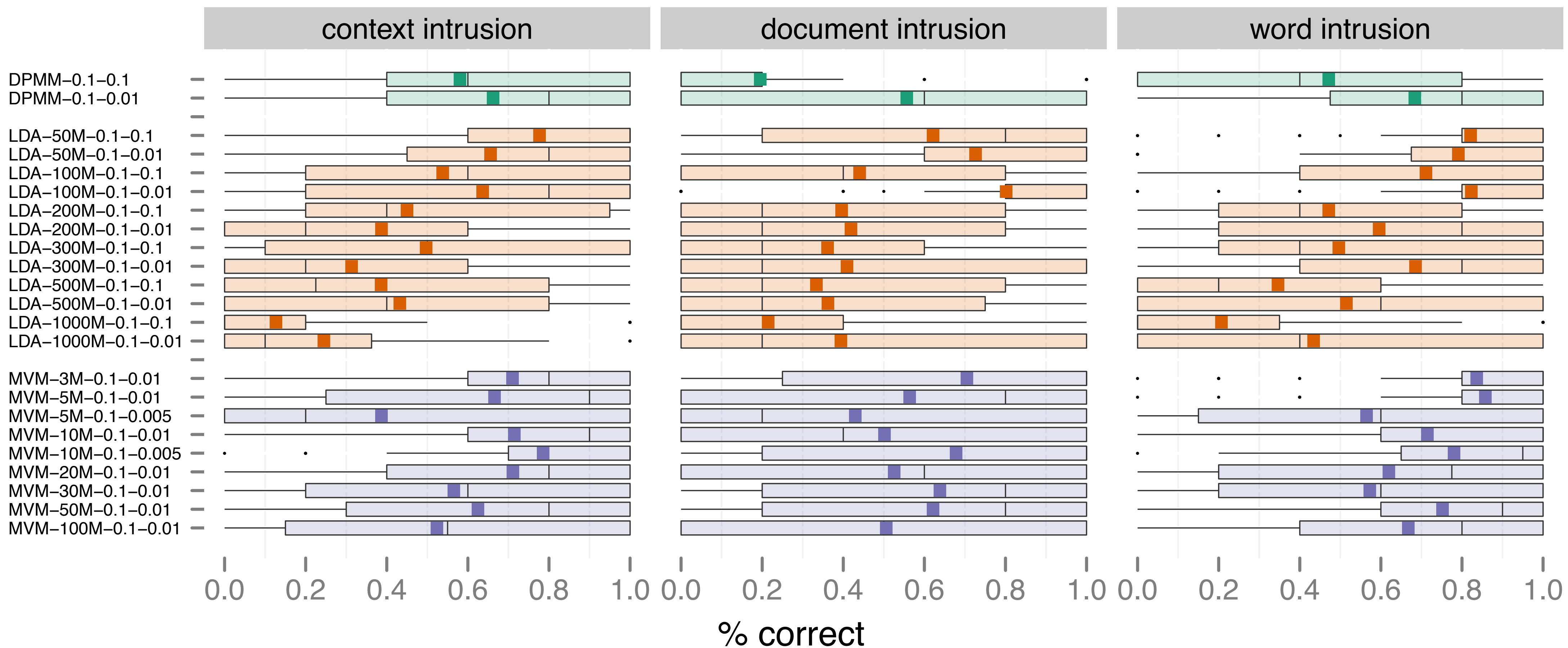
- DPMM-0.1-0.1
- DPMM-0.1-0.01
- LDA-50M-0.1-0.1
- LDA-50M-0.1-0.01
- LDA-100M-0.1-0.1
- LDA-100M-0.1-0.01
- LDA-200M-0.1-0.1
- LDA-200M-0.1-0.01
- LDA-300M-0.1-0.1
- LDA-300M-0.1-0.01
- LDA-500M-0.1-0.1
- LDA-500M-0.1-0.01
- LDA-1000M-0.1-0.1
- LDA-1000M-0.1-0.01
- MVM-3M-0.1-0.01
- MVM-5M-0.1-0.01
- MVM-5M-0.1-0.005
- MVM-10M-0.1-0.01
- MVM-10M-0.1-0.005
- MVM-20M-0.1-0.01
- MVM-30M-0.1-0.01
- MVM-50M-0.1-0.01
- MVM-100M-0.1-0.01



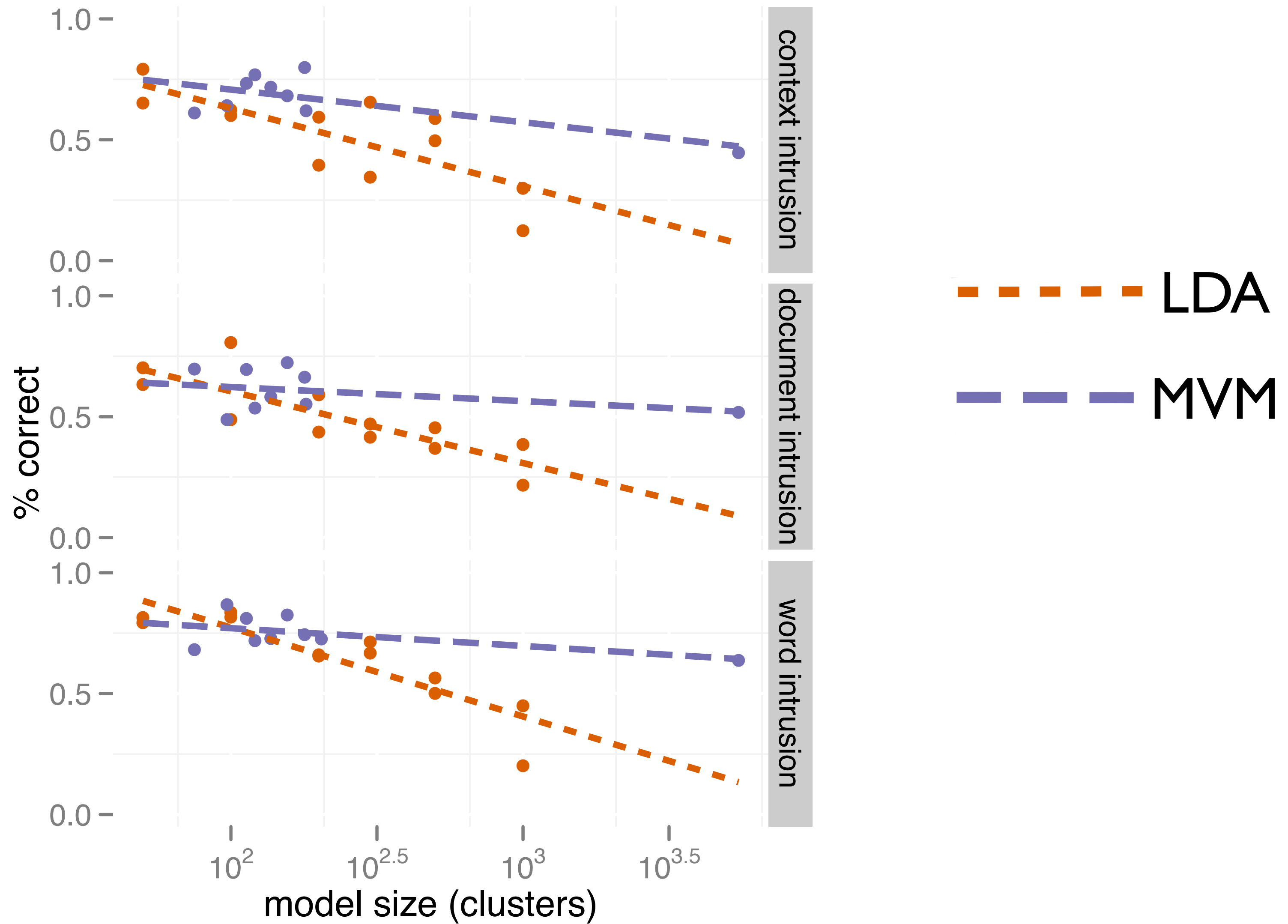
Syntax features only (freq < 50; "rare")



“Common” syntax features + document features



“Common” syntax features + document features



Conclusion

- Introduced a latent variable model accounting for cross-cutting / multiple clustering structure in word meaning
- Large-scale human evaluation of the semantic coherence of similarity predictions
- Significantly higher precision intrusion identification than related model-based approaches
- Even for fine-grained clusterings