

# Matching Consumer Health Vocabulary with Professional Medical Terms Through Concept Embedding

Yue Wang<sup>1</sup>, Jian Tang<sup>2</sup>, V.G.Vinod Vydiswaran<sup>1</sup>, Kai Zheng<sup>3</sup>, Hua Xu<sup>4</sup>, Qiaozhu Mei<sup>1</sup>

<sup>1</sup>University of Michigan; <sup>2</sup>HEC Montréal; <sup>3</sup>University of California, Irvine;

<sup>4</sup>University of Texas Health Science Center at Houston

## MOTIVATION

❖ There exists language mismatch between consumers/laypersons and health care professionals (Zeng & Tse '06). For example:

• loss of appetite  
• pain killer  
• heart attack  
• kidney stones  
• crank; ice; meth  
• ...



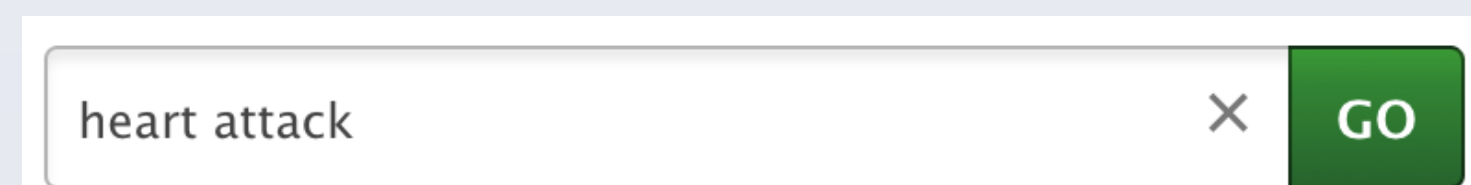
• anorexia  
• analgesic  
• myocardial infarction  
• kidney calculi  
• methamphetamine  
• ...

❖ Bridging this language gap is an important step for:

- Consumers to search for relevant health information online
- Consumers to make informed decisions
- Professionals to communicate health information & share knowledge to consumers

❖ Two previous studies on mining consumer health vocabulary (CHV):

- Consumer search queries in MedlinePlus (Zeng & Tse '06):
- Community-generated Wikipedia entries, "also known as" or "commonly known as" (Vydiswaran et al '14):



**Myocardial infarction**  
From Wikipedia, the free encyclopedia



Being active after your heart attack

Heart attack - activity; MI - activity; Myocardial infarction - activity; Car rehabilitation - activity; ACS - activity; NSTEMI - activity; Acute coronary syndrome ...  
<https://medlineplus.gov/ency/patientinstructions/000093.htm> - Medical

**Myocardial infarction (MI)**, commonly known as a **heart attack**, occurs when **blood flow** decreases or stops to a part of the **heart**, causing damage to the **heart muscle**.<sup>[1]</sup> The most common symptom

❖ These approaches obtain highly selective and accurate consumer-professional concept pairs, but many relevant pairs maybe missing, i.e. low coverage.

❖ **In this study, we aim at high coverage: given any concept in consumer health vocabulary, how can we match it to professional medical terms (and vice versa)?**

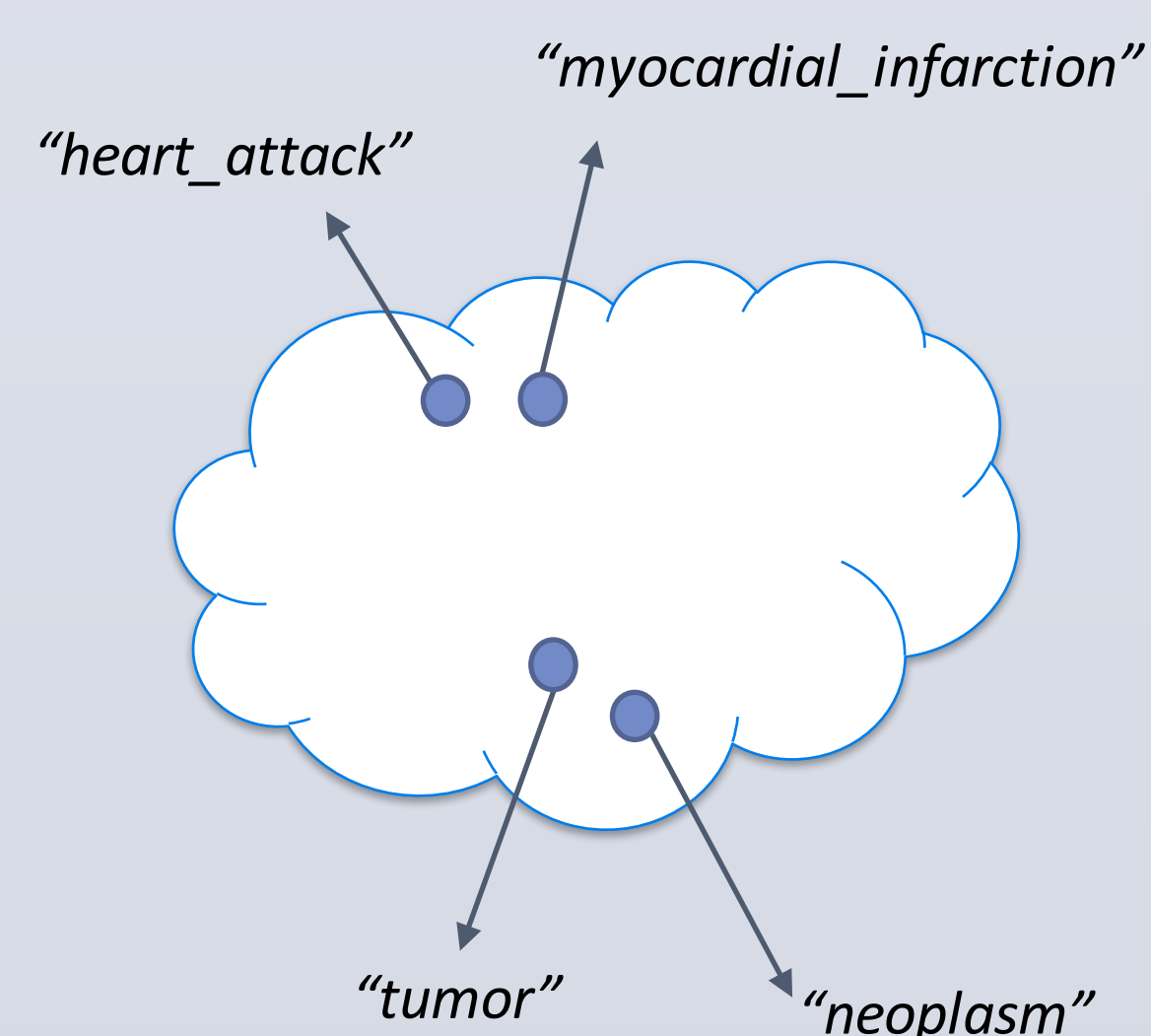
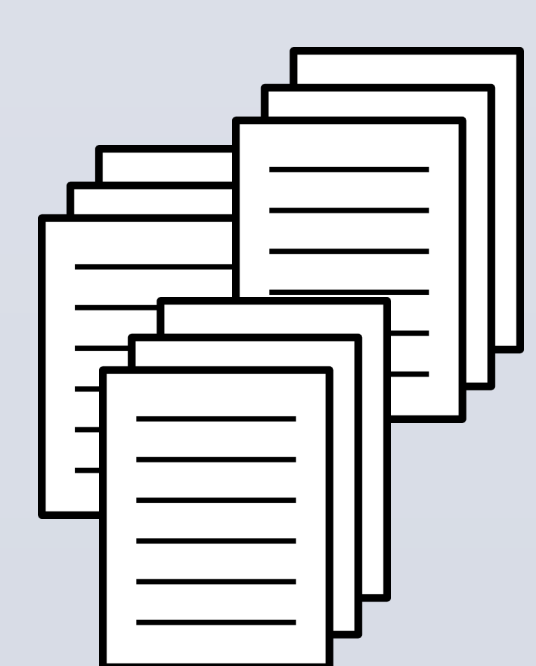
## METHODS

- ❖ We learn **vector representations** for professional and consumer medical concepts from large-scale text corpus (100 dimensions).
- ❖ Similarity between two concepts is measured by the **cosine similarity** of their vectors.

Large text collection

Word embedding learning algorithm

Concept embedding vectors



Phrase detection + Word2Vec (Mikolov '13)

❖ **Input:** three large-scale text collections:

- 2.2 million **MedHelp** posts (3 GB consumer-generated texts)
- 7 million **MEDLINE** abstracts (13 GB professional-generated texts)
- 4.7 million **Wikipedia** articles (11 GB open domain texts)

❖ **Output:** vector representation of medical concepts, including:

- 300K **words and phrases** from text collection
- 494K **medical concepts** from UMLS, including two consumer health vocabularies:
  - CHV-1 (Vydiswaran et al '14), 892 consumer/professional concept pairs
  - CHV-2 (Zeng & Tse '06), a random subset of 1,000 pairs
- Multi-word concept vector = average of word vectors

## EVALUATION & RESULTS

❖ Given a consumer concept, we return a **ranked list of concepts by decreasing cosine similarity**, and see if the corresponding professional concept is ranked high. We use Mean Reciprocal Rank (MRR) to evaluate such ranked lists.

$$\text{Reciprocal Rank (RR)} = \frac{1}{\text{hit rank}} \quad (\text{higher is better})$$

Example: Query: heart attack  
1. cardiac arrest  
2. angina  
3. myocardial infarction ✓

RR = 1/3

Table 1. Mean reciprocal rank of corresponding consumer/professional concept.

|                 | CHV-1 (Vydiswaran et al '14) |           | CHV-2 (Zeng & Tse '06) |           |
|-----------------|------------------------------|-----------|------------------------|-----------|
|                 | Con ⇒ Pro                    | Pro ⇒ Con | Con ⇒ Pro              | Pro ⇒ Con |
| Wikipedia       | 0.161                        | 0.158     | 0.403                  | 0.371     |
| MedHelp         | 0.177                        | 0.176     | 0.462                  | 0.443     |
| MEDLINE         | 0.192                        | 0.190     | 0.485                  | 0.459     |
| MedHelp+MEDLINE | 0.192                        | 0.188     | 0.483                  | 0.457     |

Con ⇒ Pro: use consumer concept as the query, compute the reciprocal rank of its corresponding professional concept. Pro ⇒ Con: vice versa.

### Example ranked lists

Query: nocturnal\_myoclonus

sleep myoclonus  
sleep disorder hypersomnia  
periodic limb movement sleep disorder  
periodic leg movements sleep  
night terrors sleep  
periodic limb movement disorder  
narcolepsy and hypersomnia  
involuntary movement extrapyramidal restlessness  
daytime somnolence sleep disorder

Query: seasonal affective disorder

seasonal affective disorder disease finding  
bipolar affective disorder  
unipolar depressive illness  
manic depressive illness  
bipolar affective disorder early onset  
bipolar disorder manic depression  
bipolar affective disorder atypical nos  
generalized anxiety  
major depressive disorder  
(expected answer: winter blues)

Query: winter blues

summer depression  
blues laundry  
winter wren  
blues music  
(expected answer: seasonable affective disorder)

### Learned word/phrase neighborhoods

Query: analgesic

analgetic  
opioid\_analgesic  
pain-relieving  
sedative  
anti-emetic  
analgesics  
non-opioid\_analgesic  
antihyperalgesic  
opioid-sparing  
narcotic\_analgesic  
analgesia  
antipyretic  
analgesic\_efficacy  
narcotic  
antitussive  
antinociceptive  
anti-nociceptive  
non-opioid\_analgesics  
sedative\_analgesic  
antiemetic

Query: painkiller

pain\_reliever  
narcotic  
painkillers  
toradol  
paracetamol  
codeine  
analgesic  
pain\_reliver  
pain\_killers  
nsaid  
muscle\_relaxer  
morphine  
tordal  
codiene  
dilaudid  
muscle\_relaxant  
pain-killer  
darvocet  
panadeine  
codine

Query: myocardial\_infarction

infarction\_mi  
acute\_myocardial  
myocardial\_infraction  
q-wave\_myocardial  
infarction\_ami  
reinfarction  
ami  
unstable\_angina  
mi  
myocardial\_infarctions  
q-wave\_mi  
infarction  
non-q-wave\_myocardial  
re-infarction  
q-wave\_infarction  
non-q\_wave  
non-q-wave\_infarction  
postmyocardial\_infarction  
postinfarction\_angina  
postinfarction

Query: heart\_attack

heartattack  
heart-attack  
heart\_attck  
heart\_attact  
stroke  
cardiac\_arrest  
angina  
heart\_attack  
heart\_atack  
heart\_attack-massive\_mi  
sudden\_death  
stoke  
sudden\_cardiac  
clogged\_arteries  
acute\_mi  
aheart\_attack  
hart\_attack  
pulmonary\_embolism  
myocardial\_infarction

❖ **Conclusion:** On average, a consumer or professional concept is 14<sup>th</sup> closest to its counterpart in CHV-1, and 7<sup>th</sup> closest to its counterpart in CHV-2.

❖ **Future work:** (1) label the concepts as consumer- or professional-oriented by propensity scores<sup>2</sup>; (2) invite medical experts for further review.

**Acknowledgments** This study is supported by NLM 2R01LM010681-05 and NSF IIS-1054199.

1. Zeng QT, Tse T. Exploring and developing consumer health vocabularies. Journal of the American Medical Informatics Association. 2006;13(1):24–29.
2. Vydiswaran VV, Mei Q, Hanauer DA, Zheng K. Mining consumer health vocabulary from community-generated text. In: AMIA. vol. 2014. American Medical Informatics Association; 2014.
3. Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J. Distributed representations of words and phrases and their compositionality. Neural Information Processing Systems, 2013.