# Growing Story Forest Online from Massive Breaking News

Presentation by **Di Niu**
University of Alberta
Computer Lab Paris 6

# Information Explosion



# Reading vs. Browsing

# News Reading: Search Engine

# News Reading: Feed Stream

## Disadvantages of existing systems

- Messed document lists

- Extremely fine-grained (articles)

- Redundant useless information

- Unstructured information

# How We Remember Information

**Event**: something revolve around one or a group of specific persons (or entities) and happen at certain place during specific time .
*Examples: Trump becomes a candidate, The first game between Kejie and AlphaGo*

**Story**: multiple events that interdependent and evolve by time form a story.
*Examples: 2016 U.S. Presidential Election, Kejie VS AlphaGo*

# How We Remember Information

**Event**: something revolve around one or a group of specific persons (or entities) and happen at certain place during specific time .
*Examples: Trump becomes a candidate, The first game between Kejie and AlphaGo*

**Story**: multiple events that interdependent and evolve by time form a story.
*Examples: 2016 U.S. Presidential Election, Kejie VS AlphaGo*

***The smallest granularity of memory: event***

# Why Event Matters

## Tags we have

**# Category tags**
**# Automotive Technology**

**# Entity tags**
**# Tesla**

## Tags we don't have

**# Event tags**
**# Tesla launches new model X**



## Title translation

Tesla: The most conscientious pricing of imported brands turned out to be it?

**7.5%** articles with event tags account for **40%** of the user traffic

# How Human Brain Organizes Information



**2016-07-19**
Trump become presidential candidate

**2016-07-26**
Hilary become presidential candidate

**2016-09-11**
Hilary attends the 911 Anniversary and leave early

**2016-09-12**
Doctor say Hilary has pneumonia

**2016-09-14**
Hilary say she was healthy

**2016-09-16**
Hilary is recovered

**2016-09-26**
First election television debate

**2016-09-28**
Hilary accuses Trump of refusing to disclose tax information

**2016-10-02**
New York Times exposures Trump tax avoidance

**2016-10-07**
Washington Post reveals Trump's speech about contempt for women

**2016-10-08**
Trump publicly apologizes for his controversial speech about women

**2016-10-10**
Second election television debate

**2016-10-19**
Third election television debate

**2016-10-28**
FBI restart "mail door" investigation

# How Human Brain Organizes Information

**2016-09-14**
Hilary say she was healthy

**2016-09-16**
Hilary is recovered

**2016-09-28**
Hilary accuses Trump of refusing to disclose tax information

**2016-10-07**
Washington Post reveals Trump's speech about contempt for women

**2016-10-08**
Trump publicly apologizes for his controversial speech about women

**2016-11-02**
Hilary condemns Trump for bullying women

**2016-11-08**
America votes to elect new president

**2016-11-09**
Donald Trump is elected president

**2016-09-26**
First election television debate

**2016-10-10**
Second election television debate

**2016-10-19**
Third election television debate

**2016-10-02**
New York Times exposures Trump tax avoidance

**2016-10-28**
FBI restart "mail door" investigation

**2016-10-29**
FBI explain for restarting "mail door" investigation

**2016-10-30**
Hilary questions FBI's motivation for restarting investigation

**2016-11-06**
FBI director: No charges after new review of Hilary emails

5  6  7  8  9  10  11  12  13  14  15  16  17  18  19  20

9

# Reinvent
# information platform
# that matches human habits

# Story Forest



Detect events automatically from massive news articles

Trees denotes stories, nodes denotes events

Edges in the tree denotes events evolving relationship

# Story Forest System Overview



**Legend:**
- w Keyword
- d Document
- e Event
- s Story

**Pipeline:**

**Preprocessing**
1. Document filtering
2. Word segmentation
3. Keyword extraction

→

**Keyword Graph**
1. Construct keyword graph
2. Community detection
3. Filtering out small sub-graphs

→

**Cluster Events**
1. Cluster by keyword sub-graphs
2. Doc-pair relation classification
3. Cluster by document graphs

→

**Cluster Stories**
1. Find the story to which each event belongs
2. Add events to existing stories, or create new stories

→

**Grow Story Forest**
1. Merge same events
2. Update story tree structure with new events

Bang Liu, Di Niu, Kunfeng Lai, Linglong Kong, Yu Xu. "Growing Story Forest Online from Massive Breaking News," in **CIKM 2017**.

# Preprocessing

Time

**Preprocessing**

1. Document filtering
2. Word segmentation
3. Keyword extraction

Table 1: Features for the classifier to extract keywords.

| Type | Features |
|---|---|
| Word feature | Named entity or not, location name or not, contains angle brackets or not. |
| Structural feature | TFIDF, whether appear in title, first occurrence position in document, average occurrence position in document, distance between first and last occurrence positions, average distance between word adjacent occurrences, percentage of sentences that contains the word, TextRank score. |
| Semantic feature | LDA |

Input Features → Gradient Boosting Decision Tree → Logistic Regression → Yes/No

# Keyword Graph



Community 1    Community 2

**Keyword Graph**

1. Construct keyword graph
2. Community detection
3. Filtering out small sub-graphs



Figure from: Yukio Ohsawa, Nels E Benson, and Masahiko Yachida. 1998. KeyGraph: Automatic indexing by co-occurrence graph based on building construction metaphor. In Research and Technology Advances in Digital Libraries, 1998. ADL 98. Proceedings. IEEE International Forum on. IEEE, 12–18.

# Cluster Events



Event 1    Event 2

**Cluster Events**

1. Cluster by keyword sub-graphs
2. Doc-pair relation classification
3. Cluster by document graphs

- Cluster by **Keyword Graph**.

- Extract **doc-pair features**: title similarity measures, content similarity measures, news category, …

- Train an **SVM classifier**: input two documents features, output if they belong to same event or not.

- Community detection on **Document Graph**

# Sentence Matching based on Deep Learning

**A. Original sentences**

Sentence A: The little Jerry is being chased by Tom in the big yard.

Sentence B: The blue cat is catching the brown mouse in the forecourt.

**B. Sentence Factorization Tree**

*AMR Purification*
*(a1)*

chase (0)
— Tom (0.0)
— Jerry (0.1)
— little (0.1.0)
— yard (0.2)
— big (0.2.0)

*Index Mapping*
*(a2)*

ROOT (0)
— chase (0.0)
— Tom (0.1)
— Jerry (0.2)
— little (0.2.1)
— yard (0.3)
— big (0.3.1)

*Node Completion*
*(a3)*

ROOT (0)
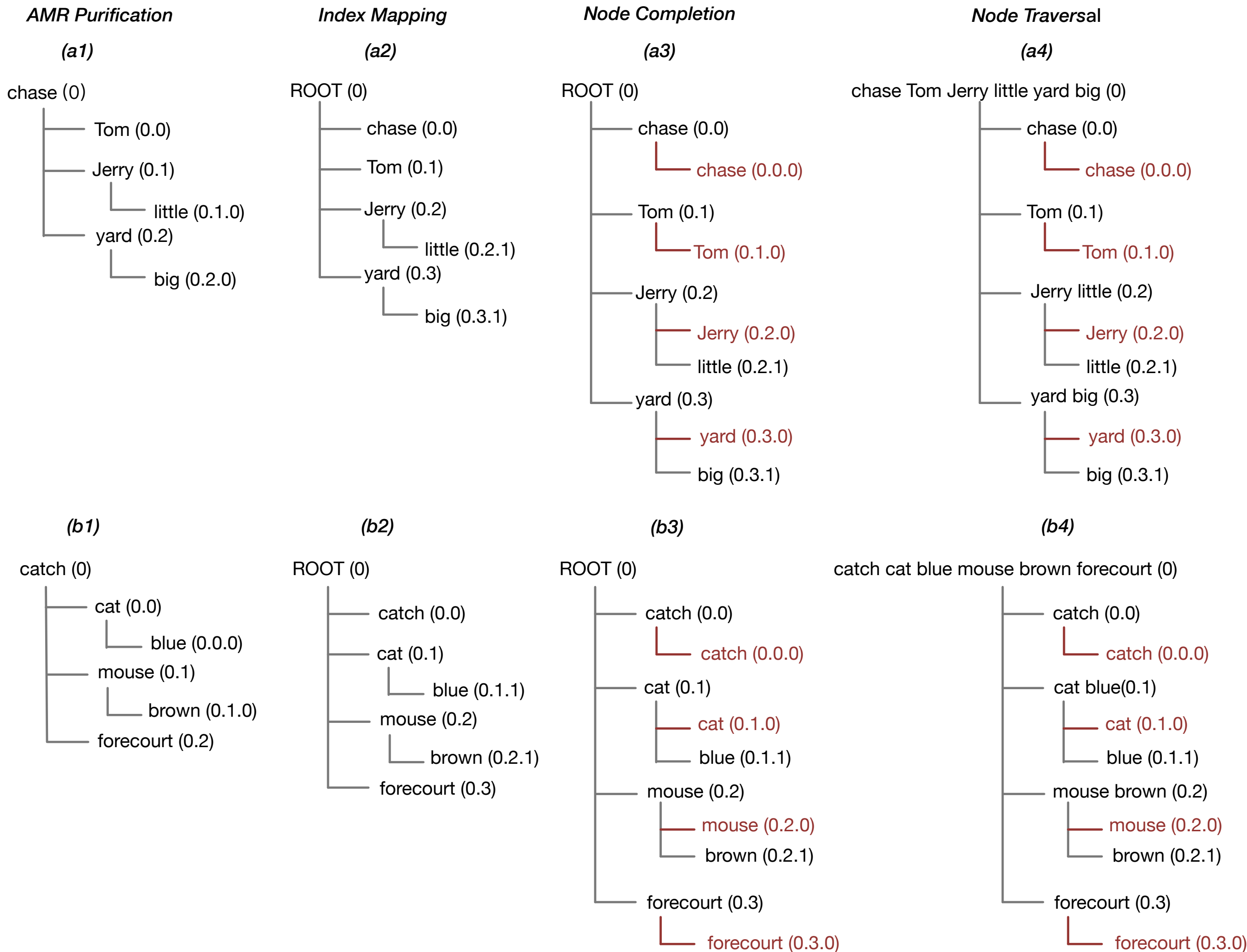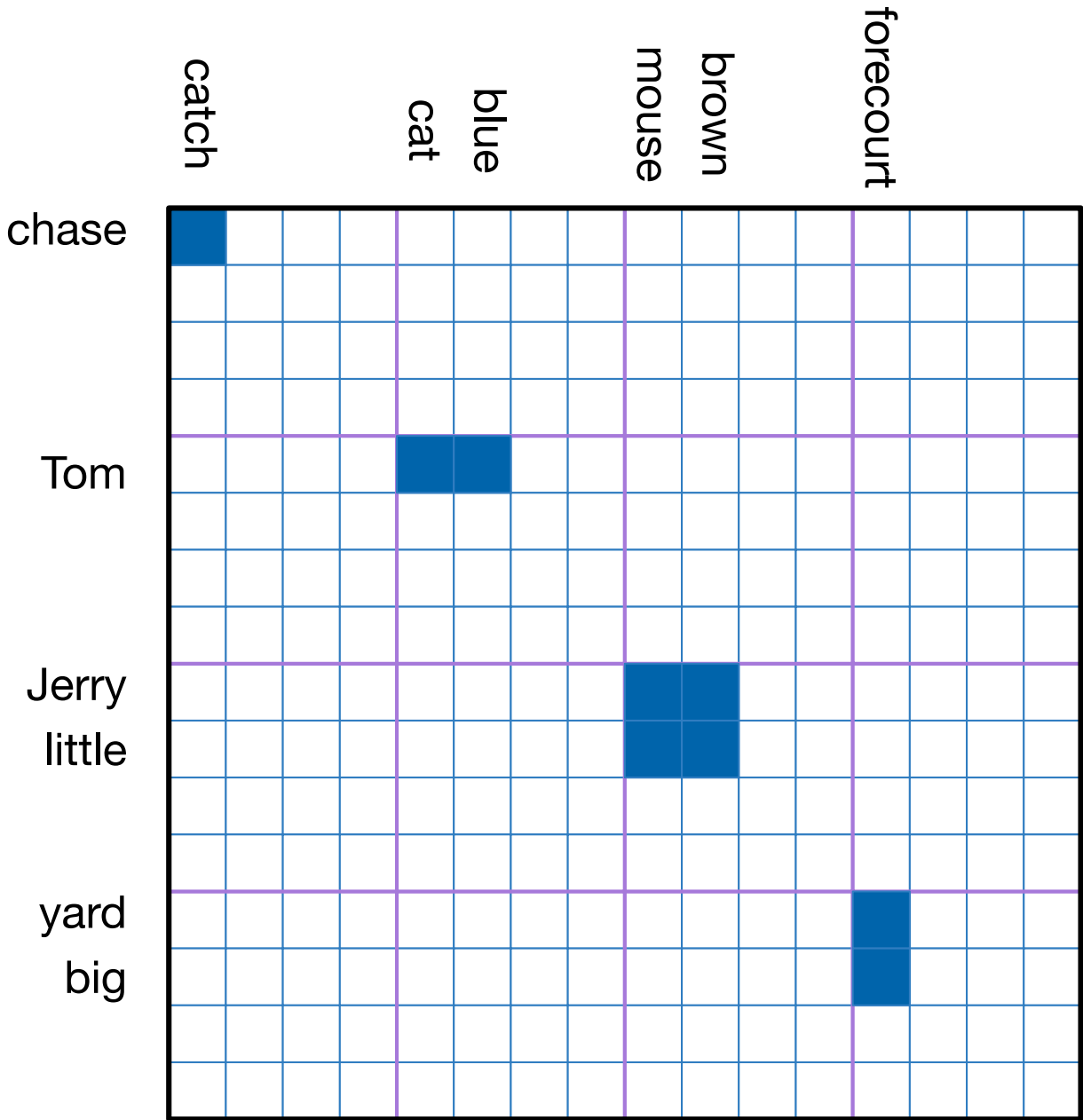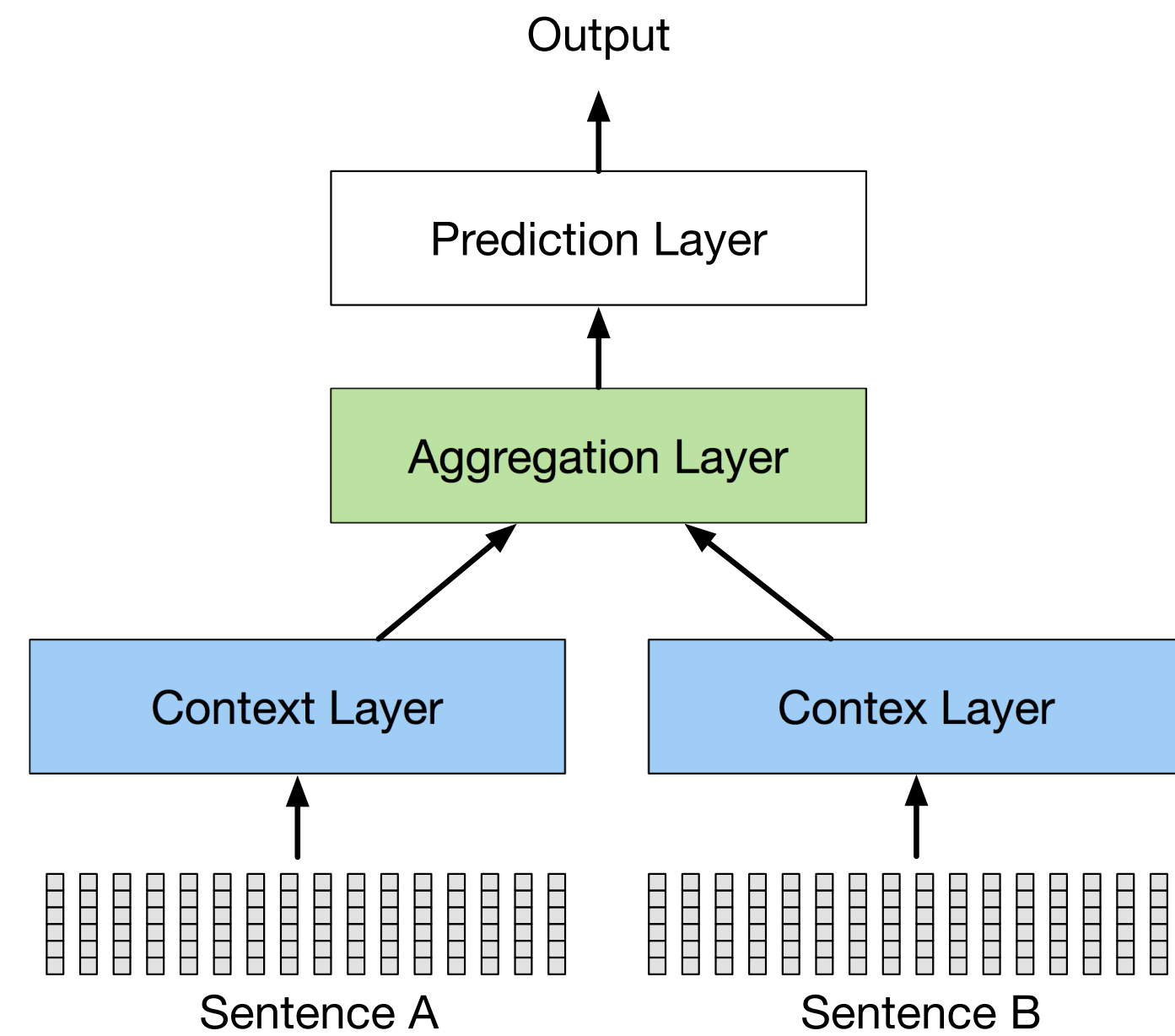— chase (0.0)
— chase (0.0.0)
— Tom (0.1)
— Tom (0.1.0)
— Jerry (0.2)
— Jerry (0.2.0)
— little (0.2.1)
— yard (0.3)
— yard (0.3.0)
— big (0.3.1)

*Node Traversal*
*(a4)*

chase Tom Jerry little yard big (0)
— chase (0.0)
— chase (0.0.0)
— Tom (0.1)
— Tom (0.1.0)
— Jerry little (0.2)
— Jerry (0.2.0)
— little (0.2.1)
— yard big (0.3)
— yard (0.3.0)
— big (0.3.1)

*(b1)*

catch (0)
— cat (0.0)
— blue (0.0.0)
— mouse (0.1)
— brown (0.1.0)
— forecourt (0.2)

*(b2)*

ROOT (0)
— catch (0.0)
— cat (0.1)
— blue (0.1.1)
— mouse (0.2)
— brown (0.2.1)
— forecourt (0.3)

*(b3)*

ROOT (0)
— catch (0.0)
— catch (0.0.0)
— cat (0.1)
— cat (0.1.0)
— blue (0.1.1)
— mouse (0.2)
— mouse (0.2.0)
— brown (0.2.1)
— forecourt (0.3)
— forecourt (0.3.0)

*(b4)*

catch cat blue mouse brown forecourt (0)
— catch (0.0)
— catch (0.0.0)
— cat blue(0.1)
— cat (0.1.0)
— blue (0.1.1)
— mouse brown (0.2)
— mouse (0.2.0)
— brown (0.2.1)
— forecourt (0.3)
— forecourt (0.3.0)

Bang Liu, Ting Zhang, Fred X. Han, Di Niu, Kunfeng Lai and Yu Xu. "Matching Natural Language Sentences with Hierarchical Sentence Factorization," in **WWW 2018.**

**D. Semantic units alignments**



16

# Sentence Matching based on Deep Learning



*(a) Siamese Architecture for Sentence Matching*    *(b) Siamese Architecture with Factorized Multi-scale Sentence Representation*

**Figure 5: Extend the Siamese network architecture for sentence matching by feeding into the multi-scale representations of sentence pairs.**

**Open Source:** https://github.com/BangLiu/SentenceMatching

Bang Liu, Ting Zhang, Fred X. Han, Di Niu, Kunfeng Lai and Yu Xu.
"Matching Natural Language Sentences with Hierarchical Sentence Factorization," in **WWW 2018.**
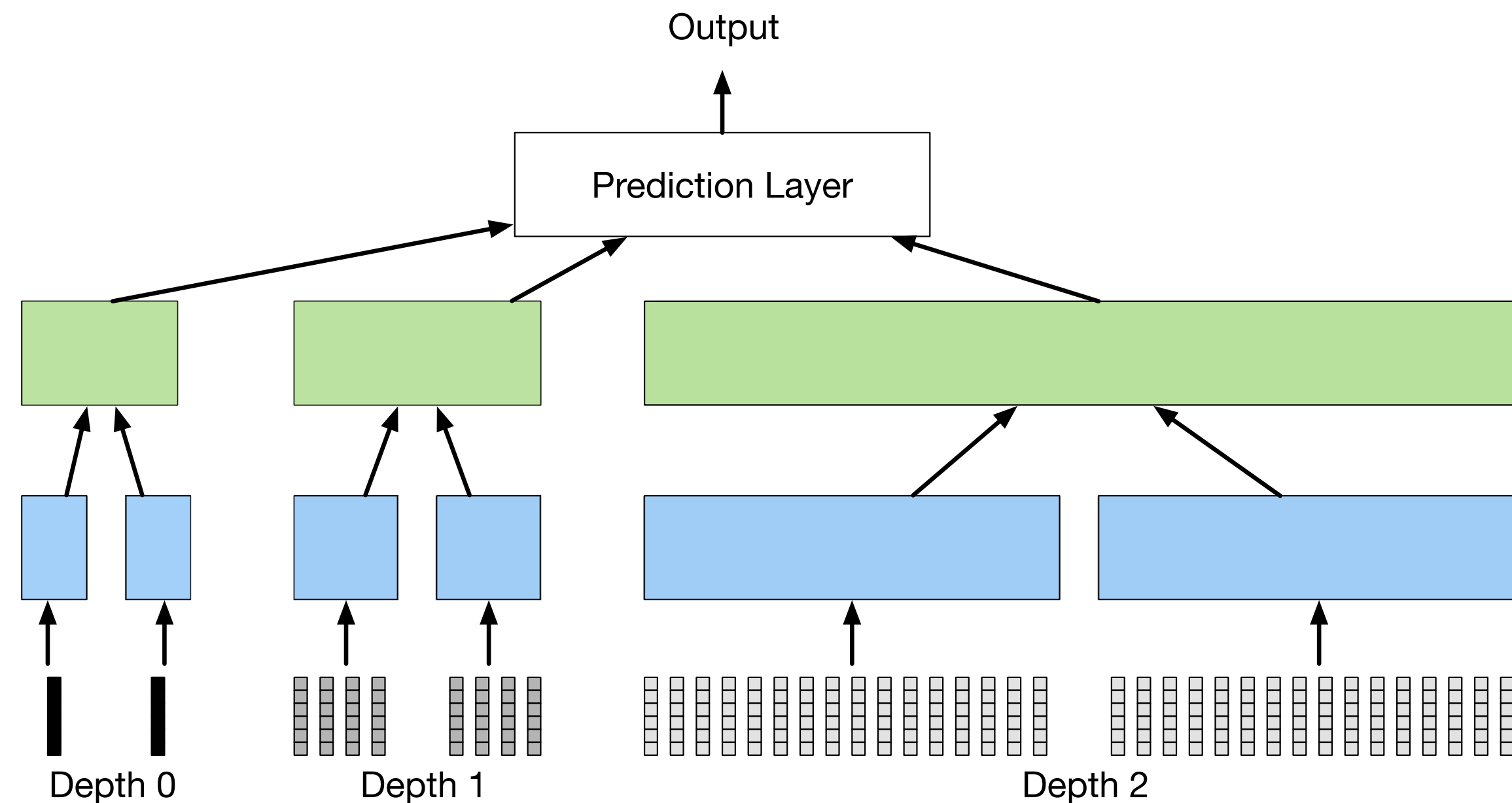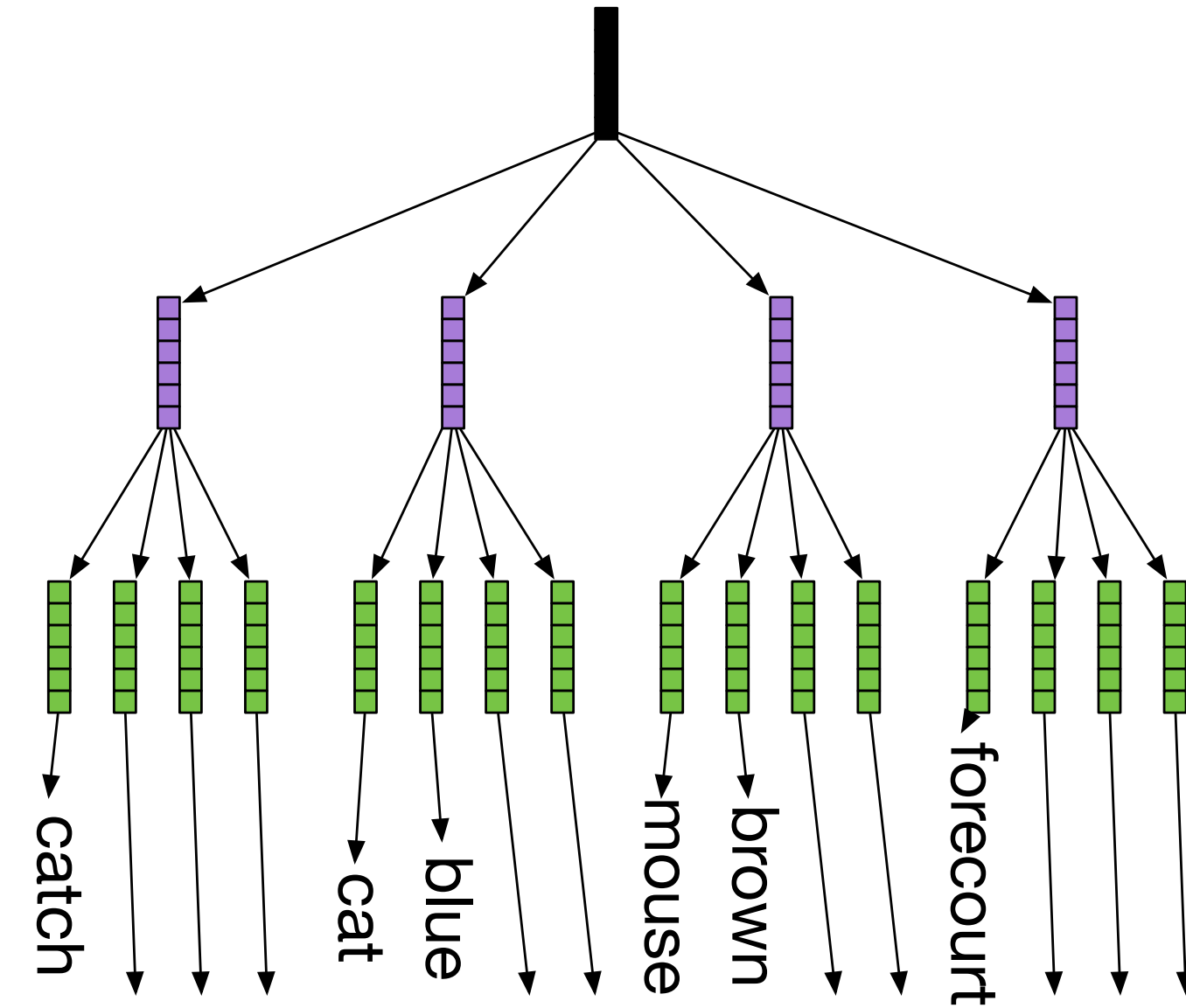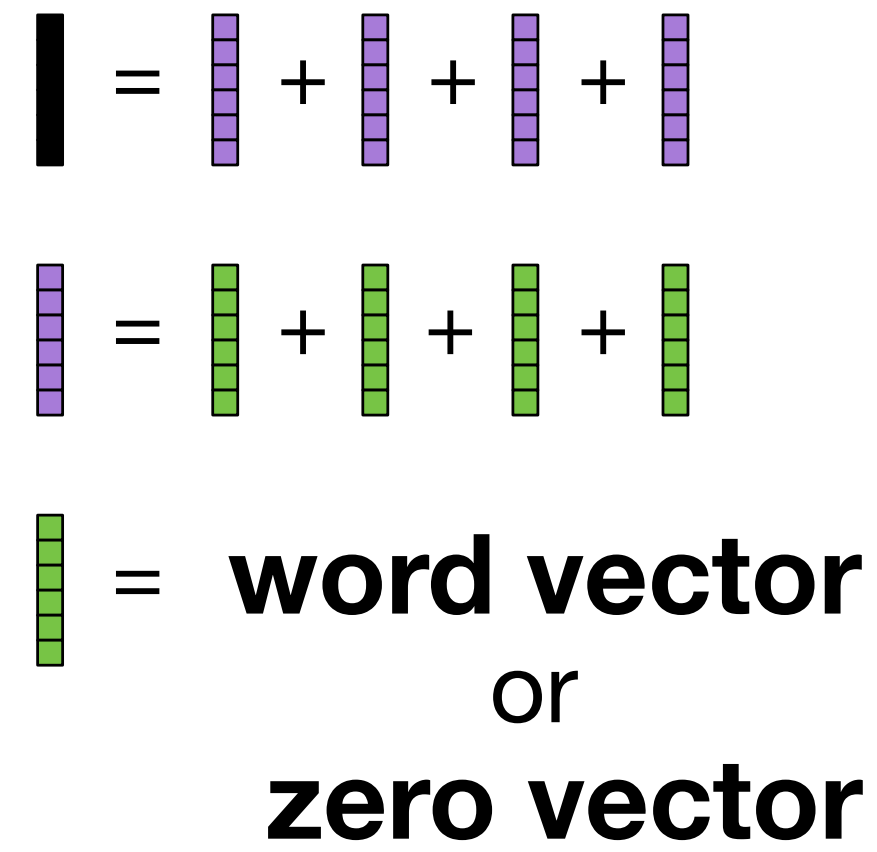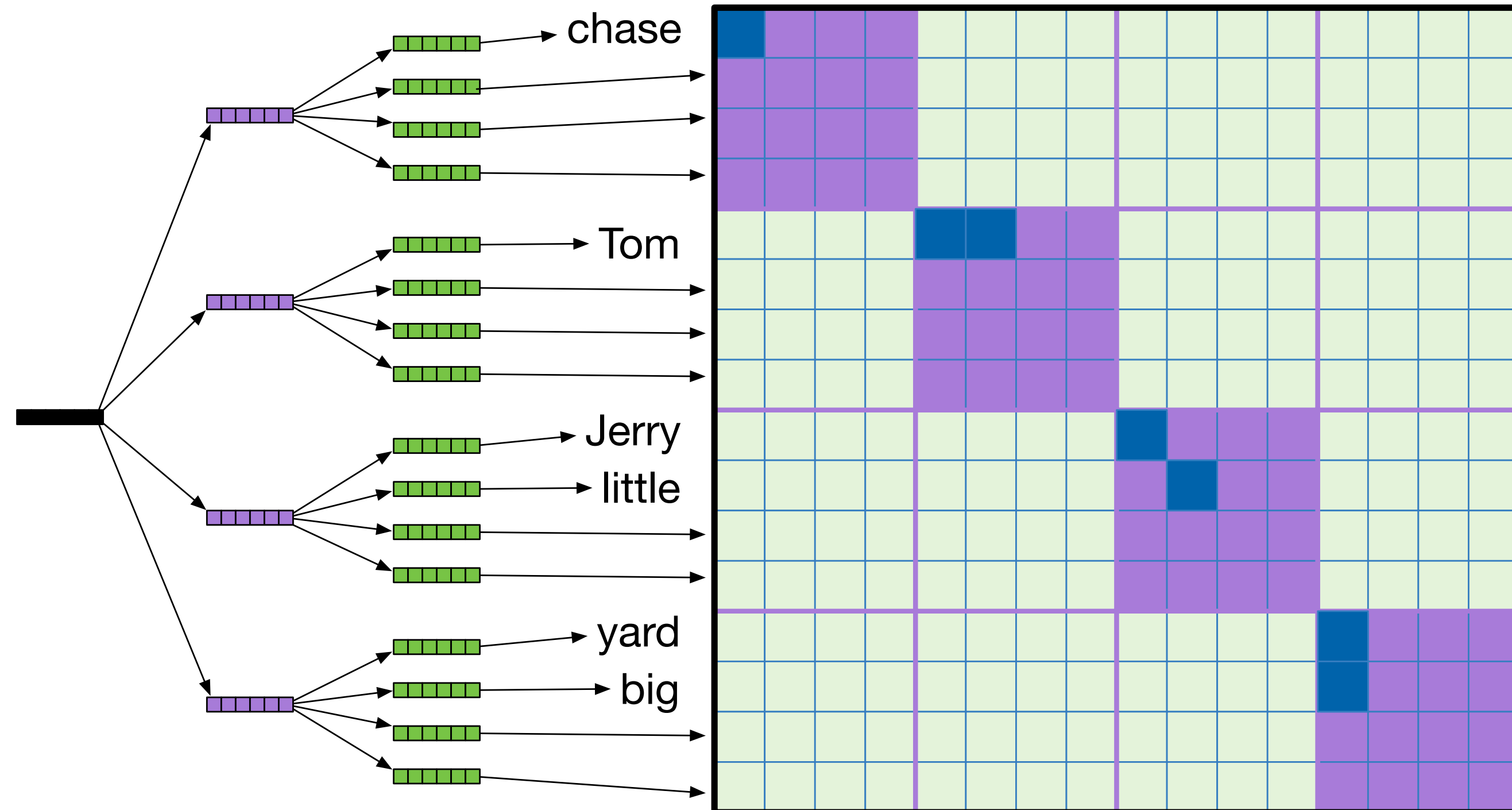
The **blue cat** is **catching** the **brown mouse** in the **forecourt**.

The **little Jerry** is **being chased** by **Tom** in the **big yard**.

catch cat blue mouse brown forecourt

chase Tom Jerry little yard big

= word vector or zero vector

Output

Prediction Layer

Depth 0    Depth 1    Depth 2

$\mathbf{I}$ = $\mathbf{I}$ + $\mathbf{I}$ + $\mathbf{I}$ +

$\mathbf{I}$ = $\mathbf{I}$ + $\mathbf{I}$ + $\mathbf{I}$ +

$\mathbf{I}$ = **word vector**
or
**zero vector**

catch    cat    blue    mouse    brown    forecourt

chase
Tom
Jerry
little
yard
big

# Long Document Matching

A graphical approach to long document matching —— **Concept Interaction Graph**

**Text:**

[1] Rick asks Morty to travel with him in the universe.
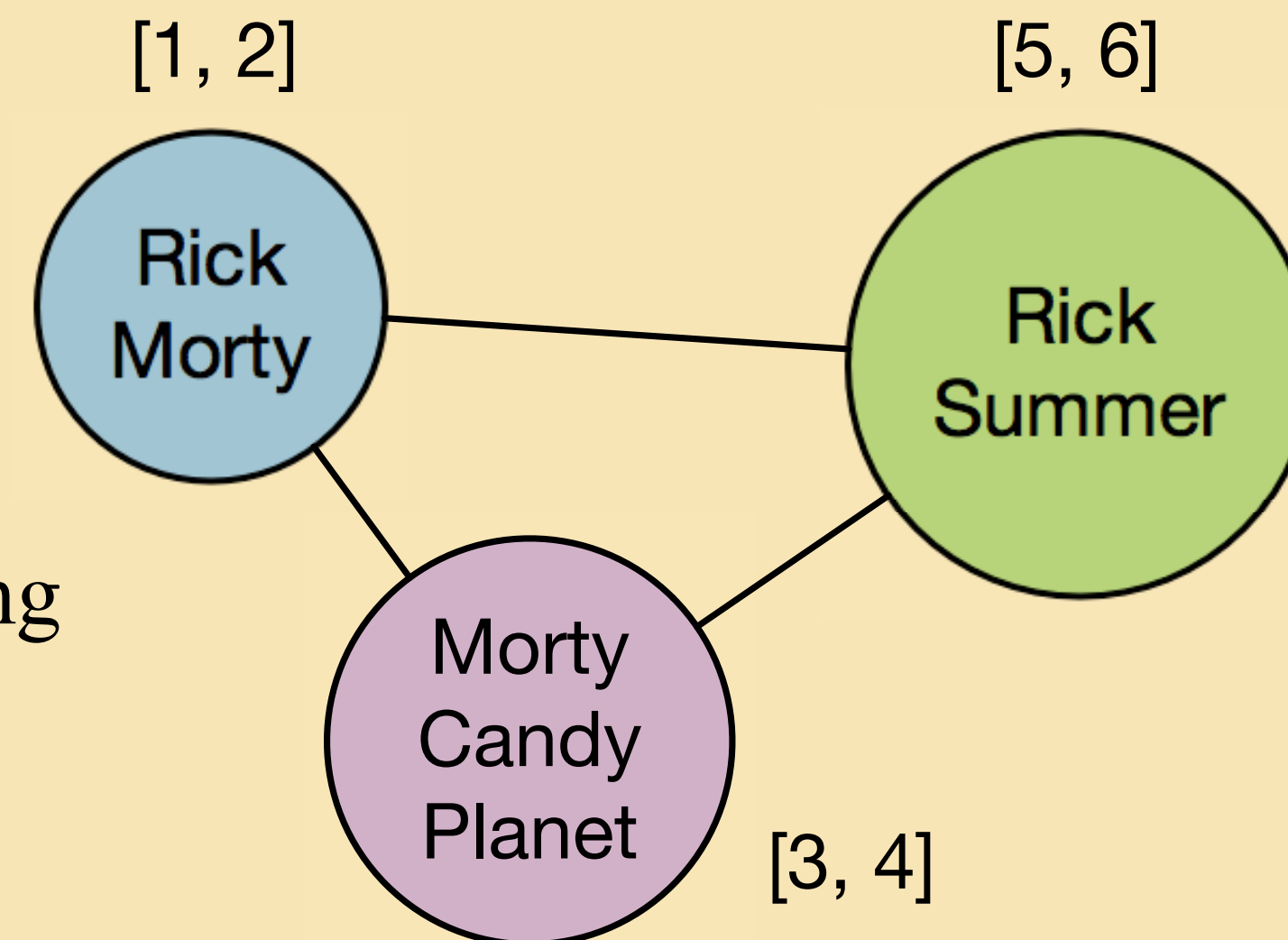[2] Morty doesn't want to go as Rick always brings him dangerous experiences.
[3] However, the destination of this journey is the Candy Planet, which is an fascinating place that attracts Morty.
[4] The planet is full of delicious candies.
[5] Summer wishes to travel with Rick.
[6] However, Rick doesn't like to travel with Summer.

**Concept Interaction Graph:**

[1, 2]       [5, 6]

Rick Morty

Rick Summer
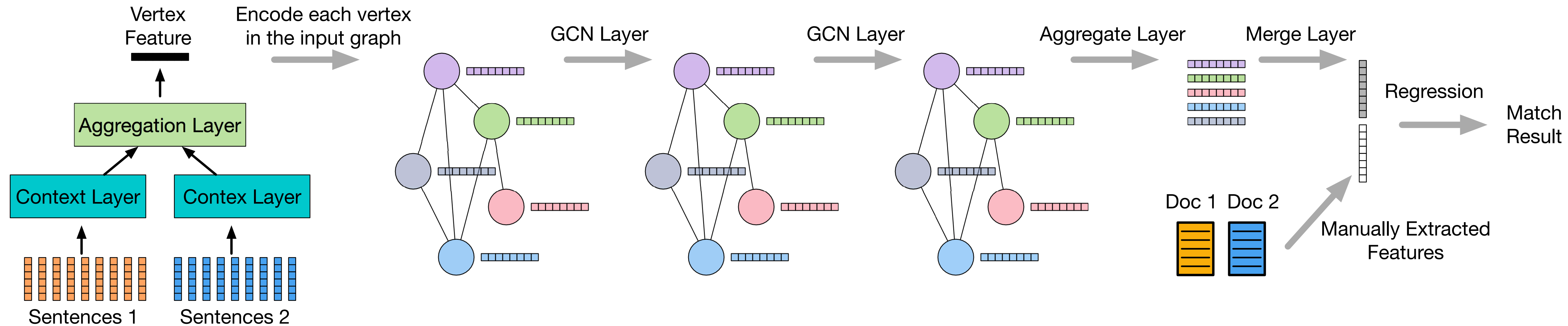
Morty Candy Planet    [3, 4]

On arXiv and under submission
**arXiv:1802.07459**

**Matching Long Text Documents via Graph Convolutional Networks**
Bang Liu, Ting Zhang, Di Niu, Jinghong Lin, Kunfeng Lai, Yu Xu

# Long Document Matching



Vertex Feature

Encode each vertex in the input graph

GCN Layer

GCN Layer

Aggregate Layer

Merge Layer

Regression

Match Result

Aggregation Layer

Context Layer

Contex Layer

Sentences 1

Sentences 2

Doc 1  Doc 2

Manually Extracted Features

*(a) Siamese Architecture for Text Pair Encoding on Each Vertex*

*(b) Architecture of Siamese Encoded Graph Convolutional Network for Long Text Pair Matching*

A graphical approach to long document matching —— **Graph Convolutional Network**

**Matching Long Text Documents via Graph Convolutional Networks**
Bang Liu, Ting Zhang, Di Niu, Jinghong Lin, Kunfeng Lai, Yu Xu

# Long Document Matching

**Table 2: Accuracy and F1-score results of different algorithms on CNSE dataset.**

| Algorithm | Dev | | Test | |
|:---:|:---:|:---:|:---:|:---:|
| | Accuracy | F1-score | Accuracy | F1-score |
| ARC-I | 0.5308 | 0.4898 | 0.5384 | 0.4868 |
| ARC-II | 0.5488 | 0.3833 | 0.5437 | 0.3677 |
| DUET | 0.5625 | 0.5237 | 0.5563 | 0.5194 |
| DSSM | 0.5837 | 0.6457 | 0.5808 | 0.6468 |
| C-DSSM | 0.5895 | 0.4741 | 0.6017 | 0.4857 |
| MatchPyramid | 0.6560 | 0.5299 | 0.6636 | 0.5401 |
| SVM | 0.7566 | 0.7299 | 0.7581 | 0.7361 |
| SE-GCN | **0.7800** | **0.7785** | **0.7901** | **0.7893** |

A graphical approach to long document matching —— **Graph Convolutional Network**

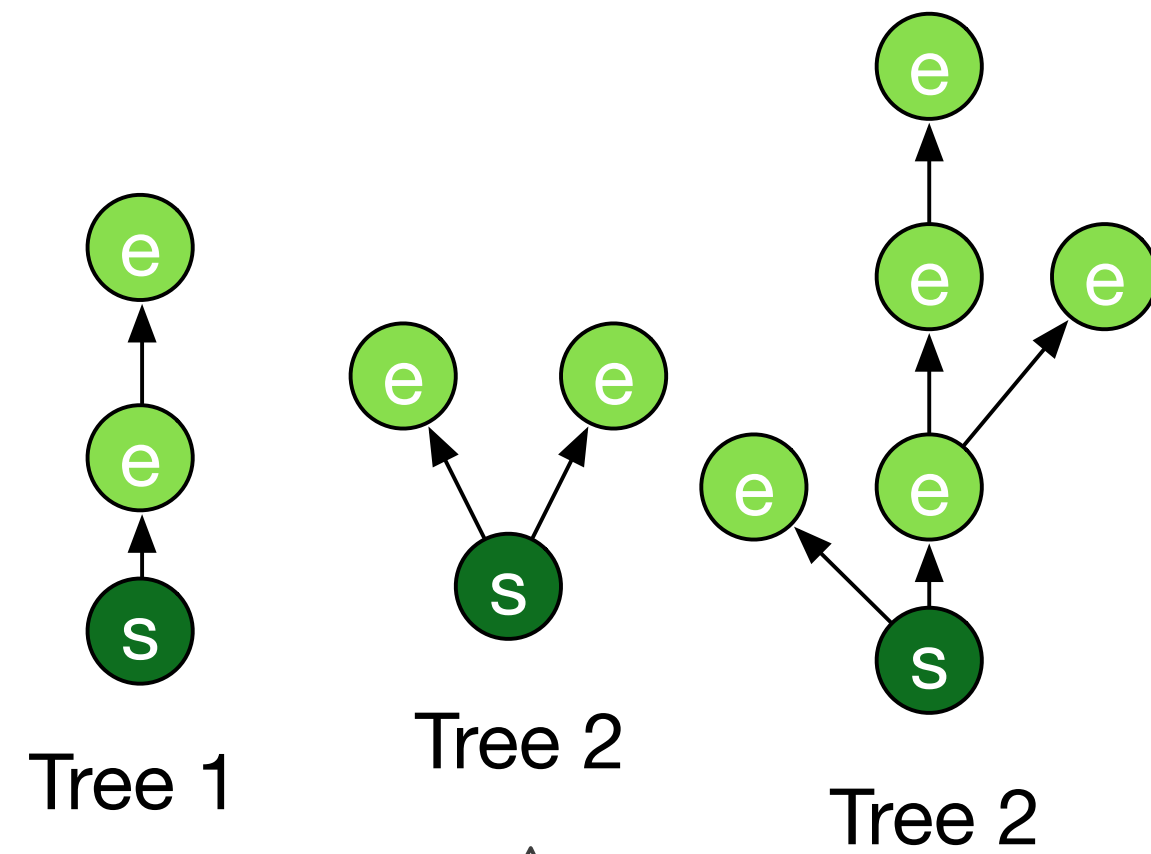**Matching Long Text Documents via Graph Convolutional Networks**
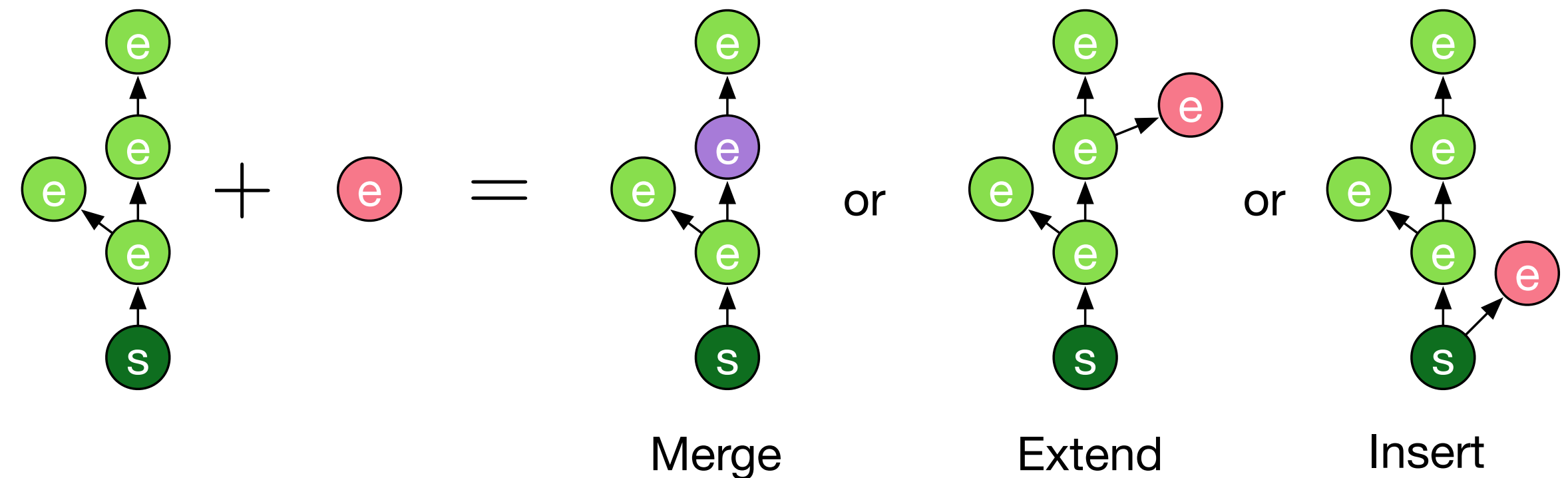Bang Liu, Ting Zhang, Di Niu, Jinghong Lin, Kunfeng Lai, Yu Xu

# Cluster Stories

Story 1

Story 2

Cluster Stories

1. Find the story to which each event belongs
2. Add events to existing stories, or create new stories

**Story**: multiple events that are interdependent and evolve over time form a story.

# Story Structure Generation



Tree 1

Tree 2

Tree 2

## Grow Story Forest

1. Merge same events
2. Update story tree structure with new events

Event    New event    Merged event    Story

Merge      Extend      Insert

**Choose the best position in the tree to insert a new event node**

# Clustering Performance

- **LDA+Affinity Propagation**: extract 1000 dimensional LDA feature, clustering by Affinity Propagation.

- **KeyGraph**: the original KeyGraph algorithm proposed in [1], which doesn't include the second step in our approach.
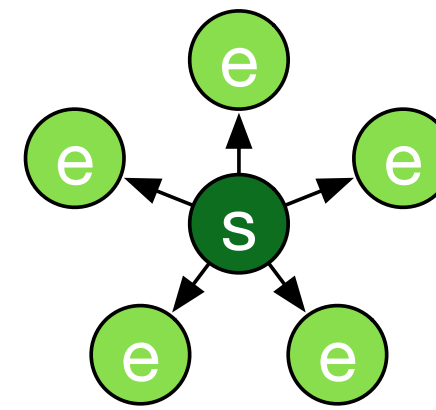
- **StoryForest**: our approach.

**Table 2: Comparing different event clustering methods.**

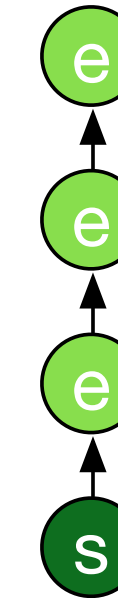| Algorithm | Homogeneity | Completeness | V-measure |
|---|---|---|---|
| Our approach | **0.960** | 0.965 | **0.962** |
| KeyGraph | 0.554 | **0.989** | 0.710 |
| LDA + AP | 0.620 | 0.947 | 0.749 |

Bang Liu, Di Niu, Kunfeng Lai, Linglong Kong, Yu Xu. "Growing Story Forest Online from Massive Breaking News," in **CIKM 2017**.
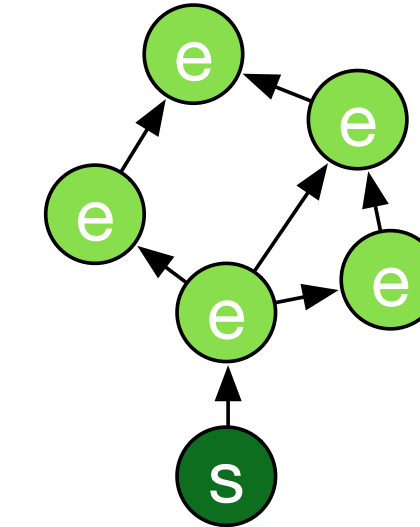
# Story Structure Performance

- **Flat Cluster:** cluster by stories, no structure.

- **Story Timeline:** organizes events linearly by time.

- **Story Graph:** calculates a connection strength for each pair of events and connect the pair if the score exceeds a threshold.

- **Event Threading:** appends each event to its most similar earlier event. Similarity measured by TF-IDF.
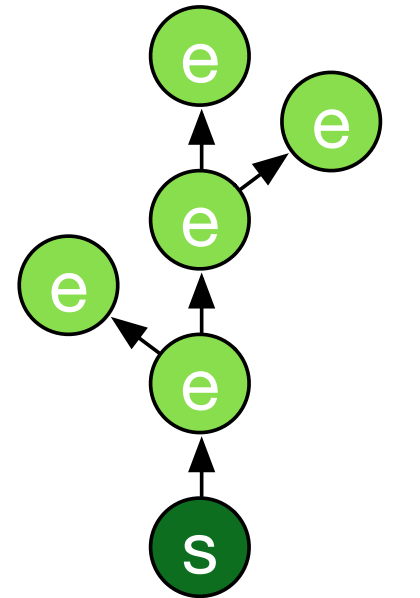


(a) Flat structure   (b) Timeline structure   (c) Graph structure   (d) Tree structure

**Table 3: Comparing different story structure generation algorithms.**

|  | Tree | Flat | Thread | Timeline | Graph |
|---|---|---|---|---|---|
| Correct edges | **82.8%** | 73.7% | 66.8% | 58.3% | 32.9% |
| Consistent paths | **77.4%** | – | 50.1% | 29.9% | – |
| Best structure | **187** | 88 | 84 | 52 | 19 |

(from the **CIKM 2017** paper)

26

# Deployed in Tencent QQ browser
## The hot topic list

# Dr. Hawking's PhD thesis made public

---

**中国联通** 11:10 AM

🔍 热点 ✕ 搜索热点　　取消

热搜榜

1. 女子坐飞机唯一乘客
2. 楼市出"王炸"
3. C罗蝉联足球先生
4. 逛菜市怕弄脏萨摩
5. 左右脑年龄测试不靠谱
6. 女子带宝宝自考
7. 霍金公开博士论文
8. 迪拜警察新座驾
9. 90后毕业写小说
10. 蒂勒森突访阿富汗
11. 6岁娃娃独自撑起一个家

---

**中国联通** 8:46 AM

**#霍金公开博士论文#**
1054阅 5人参与
主持人：永恒代價 ›
关注

导语：你是不是也很好奇，在史蒂芬·霍金还未成为博士的时候，他是如何看待黑洞理论的？ 现在，这个问题有答案了。

澎湃新闻
10月24日 2阅

精华 霍金首次公开24岁时博士论文，希望激发大众对星空探索兴趣

你是不是也很好奇，在史蒂芬·霍金还未成为博士的时候，他是如何看待黑洞理论的？ 现在，这个问题有答案…

分享　　我来说两句　　👍 2

---

**中国联通** 8:47 AM

界面新闻
10月23日 0阅

精华 霍金首次免费公开博士论文 勉励全球人民"仰望星空"

今年3月，史蒂芬·霍金因其在理论物理和宇宙学方面的卓越成就获"伦敦市自由荣誉市民"奖。图片来源：东方…

分享　　评论　　赞

孙若空
10月24日 0阅

精华 「这世界」剑桥大学首次公开了霍金的博士论文，去看吗？

昨天，物理学家史蒂芬·霍金首次公开了自己 1966 年写的博士论文。零点刚过，剑桥大学的网站就上线了霍金这…

我来说两句

# Thanks!
## Q&A