



应用：语言自动生成

一、自然语言生成概述

自然语言生成定义

自然语言生成历史

自然语言生成数据集

自然语言生成现有系统

自然语言生成应用

NLG(NATURAL LANGUAGE GENERATION)

维基百科定义

- **Natural language generation (NLG)** is the natural language processing task of generating natural language from a machine representation system such as a knowledge base or a logical form

按照不同的输入划分

- 文本到文本的生成(text-to-text generation)
- 意义到文本的生成(meaning-to-text generation)
- 数据到文本的生成(data-to-text generation)
- 图像到文本的生成(image-to-text generation)

NLG 发展历史

模板生成技术

模式生成技术

短语规则扩展技术

属性特征生成技术

模板生成技术

最早采用的自然语言生成技术

设计可能出现的语言情况，构造相应模板(包括常量、变量)

根据用户输入信息替代模板中变量，生成文本

优点：效率高，实现手段简单

缺点：处理仅在字符级上处理，生成文本质量不高，难以满足多变的需求

模式生成技术

模式生成技术(**Schema based generation**)是基于语义学中的修辞谓词来表达文本结构

- 文本表示成结构树形式(**Root, Schema, Predicate, Argument, Modifier**)
- **Root**是根节点，表示一篇文章
- **Schema**是子节点，表示一段话或几句话
- **Predicate**是一个子树，表示一个句子(文章的基本单位)
- **Argument**是叶子节点，表示句子中的基本语义成分
- **Modifier**是叶子节点，代表具有对修饰成分**Argument**

具有较好的维护性，生成的文本质量高，但只能用于固定结构段落的生成

短语规则扩展技术

基于结构修辞理论(RST)，文章的各个组成部分都是由一些特定的关系按照一定的层次内聚在一起

包含两种模式：**nucleus-satellite**和**multi-nucleus**

- **nucleus-satellite**包括表达基本命题和表达附属命题，其组合表达目的、因果、转折、背景等关系
- **multi-nucleus**涉及一个或多个语段，用于说明顺序、并列等关系

具有更强的灵活性，在生成文本时也生成文本的总体框架结构，缺点在于基本数据结构、文本规则库较难建立

属性特征生成技术

在生成系统中，每个变化都是由一个属性特性表示出来

- 生成文本是主动or被动
- 生成的文本表示的动作是问题or命令

输出单元与特定的属性特征集相连，在生成过程中对每个信息增加对应的属性特征，确定输出结果

- 特征属性一般是语法特征
- 输出单元是词汇

优点在于概念简单，生成文本灵活，但难以维护各个属性间的内容关系，难以控制特征集选择

NLG 数据集

Generation Target

- **PIL: Patient Information Leaflet corpus**
 - 可搜索可浏览的具有SGML标注的各种文档格式的患者信息传单

Content selection, aggregation

- **SumTime Meteo**
 - 天气预报文本和基于的数值数据平行语料库
- **Wizard-of-Oz corpus**
 - 爱丁堡订餐领域的对话系统语料库

NLG 数据集

Generating referring expressions

- **COCONUT Corpus**
 - 在购买家具过程中的对话系统中的语料构成的数据集
- **GRE3D3 and GRE3D7: Spatial Relations in Referring Expressions**
 - 分别包含**720**个和**4480**个指称表达式，每个指称表达式是描述一个对象在三维空间的场景
- **TUNA Reference Corpus**
 - 语义和语用上的指称表达式，用于识别视觉上的对象的语料库

NLG现有系统

基于Solog的深层自然语言生成器系统—[ASTROGEN](#)

混合模板/基于Word的生成系统(生成商务信函)—[CLINT](#)

应用SegSim策略生成英语句子--[NLGen](#)

基于n-gram和hypertagging, 用于对话系统—[OpenCCG](#)

基于词汇化的语法构建句子的语法, 可作句子规划和生成——[SPUD](#)

可以为临床医学文件生成文本的系统——[Suregen-2](#)

NLG应用

论文写作

摘要生成

自动作诗

新闻写作

报告生成

百科写作

.....

论文自动生成

SCIgen- An Automatic CS Paper Generator

只要输入作者名，就可以生成“SCI级别”的computer science论文

SCIgen多次成功逆袭IEEE的国际会议

- 在2005年，机器论文Router: A Methodology for the Typical Unification of Access Points and Redundancy被WMSCI会议所接收
- 2008年和2009年中国武汉举办的两个IEEE国际会议投稿，还获得高度评价



摘要生成

新闻摘要应用 Summly

- 提取新闻摘要功能，可以通过书签工具把文章发送到 Summly
- 2013年雅虎收购新闻摘要应用 Summly

微软“万小冰”提供金融领域公告摘要服务

- 实时抓取沪深两市公告作为基础数据
- 处理文件中表格、文字、数据等
- 对公告内容进行分类和结构化处理
- 建立自动摘要生成模型
- 在公告发布后**20秒**左右生成高质量摘要

科大国创:2018年半年度业绩预告

公告日期: 2018-07-14

发布时间: 2018-07-13 19:32

附件: [科大国创: 2018年半年度业绩预告.pdf](#)

公告摘要:

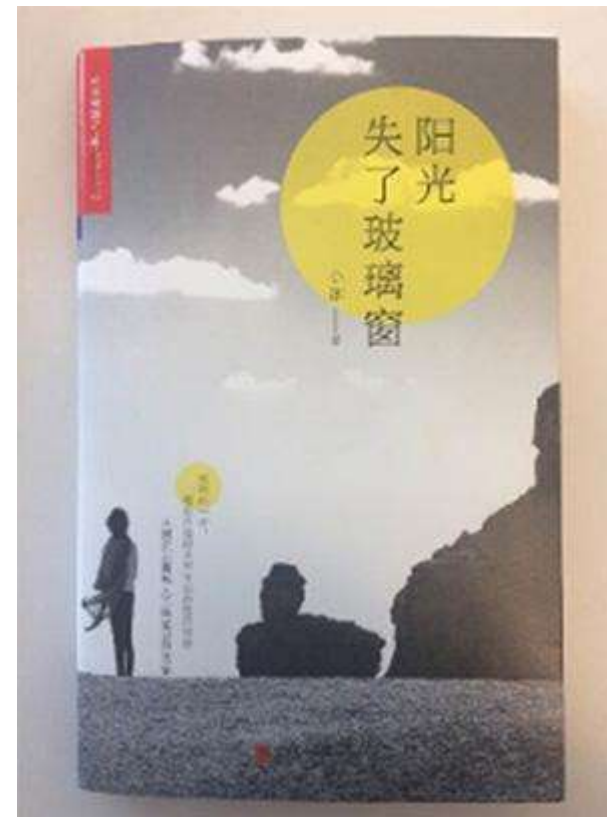
科大国创2018年半年度业绩预告,预计归属于上市公司股东的净利润,比上年同期增长25%—55%、盈利581.14万元—720.61万元。主要原因是 因公司业务存在明显的季节性波动,通常来说,公司营业收入及相应销售回款上半年较少,但期间的相关费用并没有减少,从而导致公司净利润的季节性波动明显。报告期内,公司整体经营平稳。公司预计报告期内非经常性损益对归属于上市公司股东净利润的影响金额约为300万元左右。(摘要来自万小冰)

自动作诗

清华大学作诗机器人“薇薇”通过“图灵测试”

- “薇薇”创作的诗词中，有**31%**被认为是人创作
- “薇薇”通过社科院等唐诗专家评定，通过“图灵测试”

聊天机器人小冰创作的诗集《阳光失了玻璃窗》出版



新闻写作

Automated Insights公司开发的wordsmith平台可以在每秒之内生成近2000篇新闻稿件

2015年，腾讯通过机器人新闻写作可以在政府发布CPI资料之后，只用了几分钟的时间就完成了相关新闻稿件的发布

今日头条和南方都市报先后与北大计算机所合作，分别推出奥运AI小记者Xiaomingbot和“小南”写稿机器人



报告生成

「心声医疗」帮医生自动生成诊断报告

- 通过CNN、RNN等AI技术分析心电图，生成诊断报告

CMU 邢波教授利用 AI 自动生成医学影像报告

- 运用图像说明技术(CNN-RNN框架)，可以为胸部X射线影像添加文本标签,自动生成文字描述



百科生成

2009年提出结构化方法生成维基百科文章

- 采用产生相关主题的模板，根据主题选择内容构成百科

谷歌大脑提出通过多文档摘要方法生成维基百科

- 输入维基百科的主题和参考文献的集合，目标是生成维基百科文章的文本
- 看作多文档摘要任务，并应用了**Transformer** 结构

小结

自然语言生成概述

- 定义及划分的四种类型
- 发展历史阶段
- 数据集
- 现有的系统
- 应用场景及示例

二、数据到文本的生成技术

数据到文本生成的定义

数据到文本生成系统框架

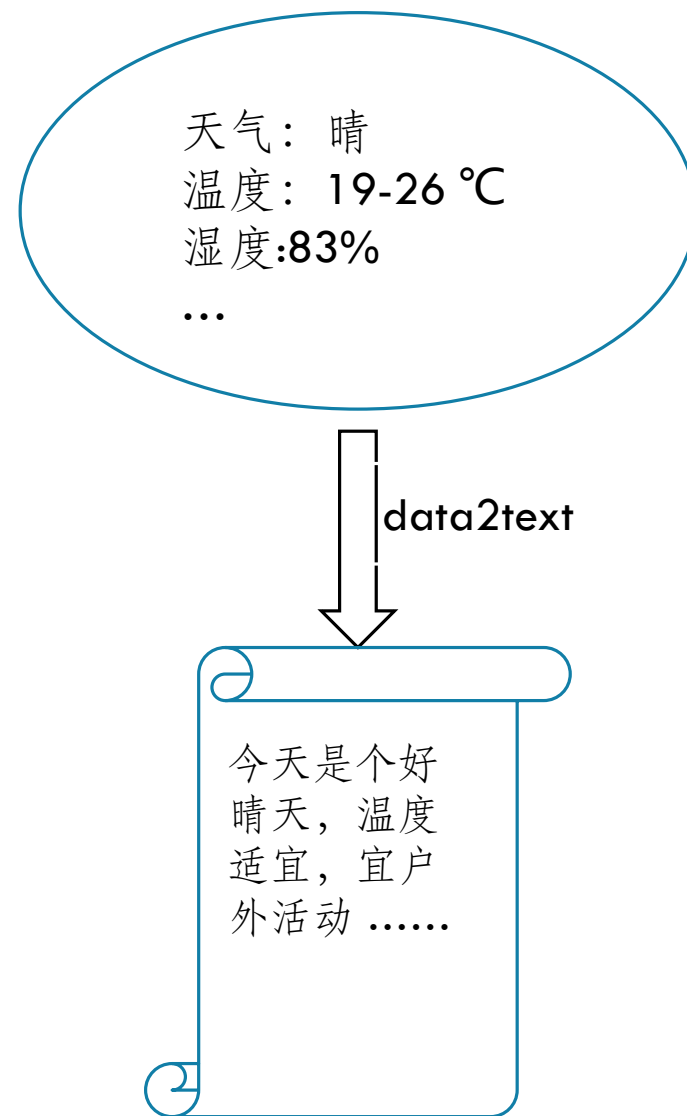
数据到文本生成的应用

数据到文本生成的研究前沿

小结

定义

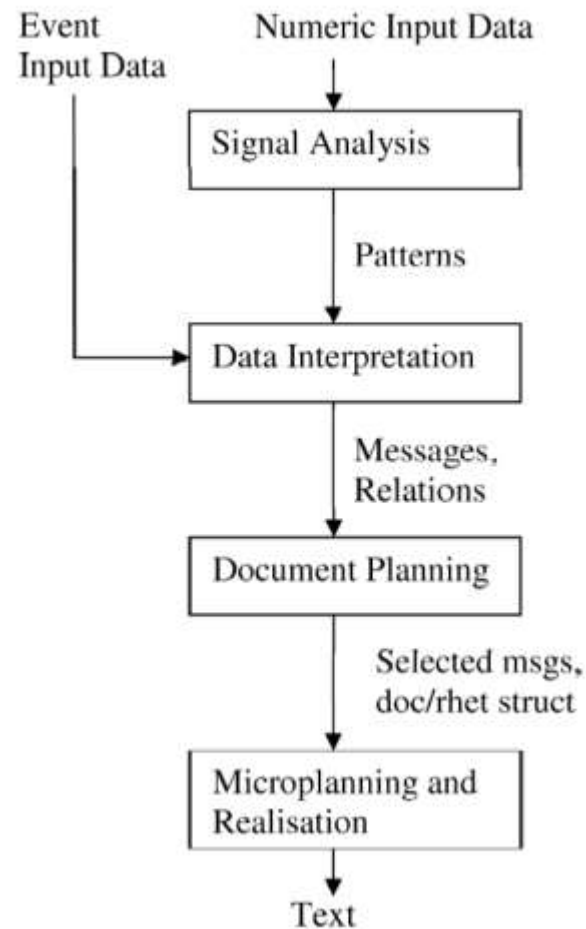
数据到文本的生成技术指根据给定的数值数据生成相关文本，例如基于数值数据生成天气预报文本、体育新闻、财经报道、医疗报告等



数据到文本生成系统框架

英国阿伯丁大学的 **Ehud Reiter** 在三阶段流水线模型：

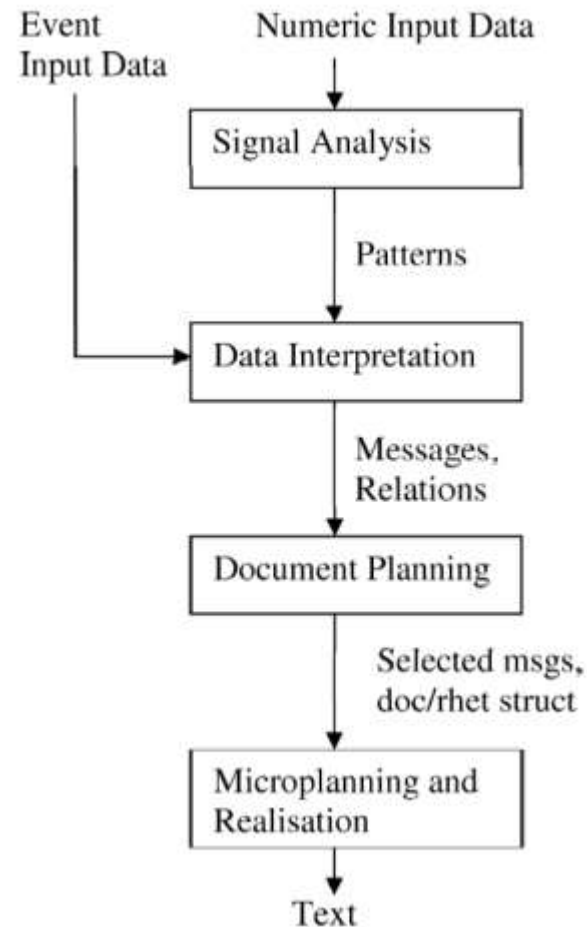
- 信号分析模块(Signal Analysis)
- 数据阐释模块(Data Interpretation)
- 文档规划模块(Document Planning)
- 微规划与实现模块(Microplanning and Realisation)



数据到文本生成系统框架

■ 信号分析模块(Signal Analysis)

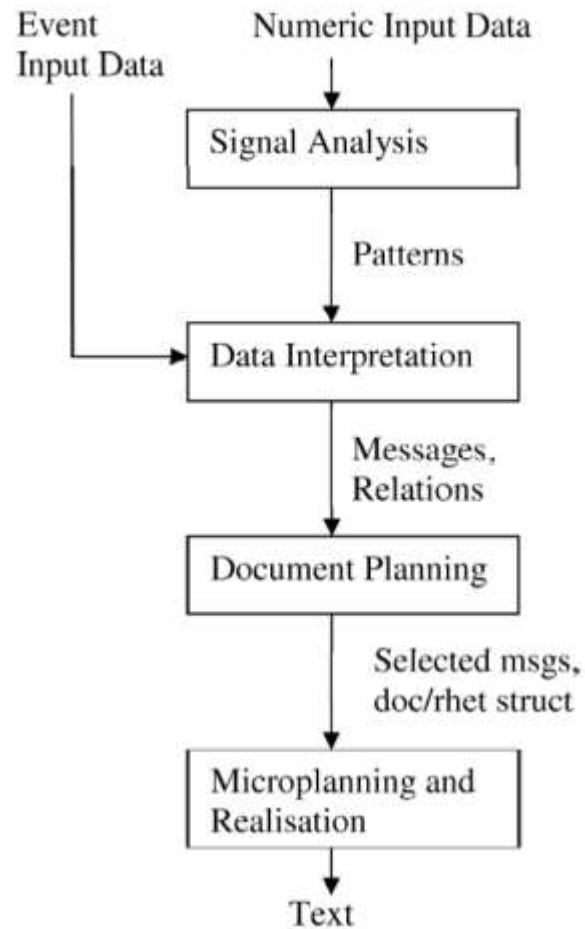
- 输入为数值数据，通过利用各种数据分析方法检测数据中的基本模式，输出离散数据模式
- 该模块与具体应用领域和数据类型相关，针对不同的应用领域与数据类型所输出的数据模式是不同的



数据到文本生成系统框架

■ 数据阐释模块(Data Interpretation)

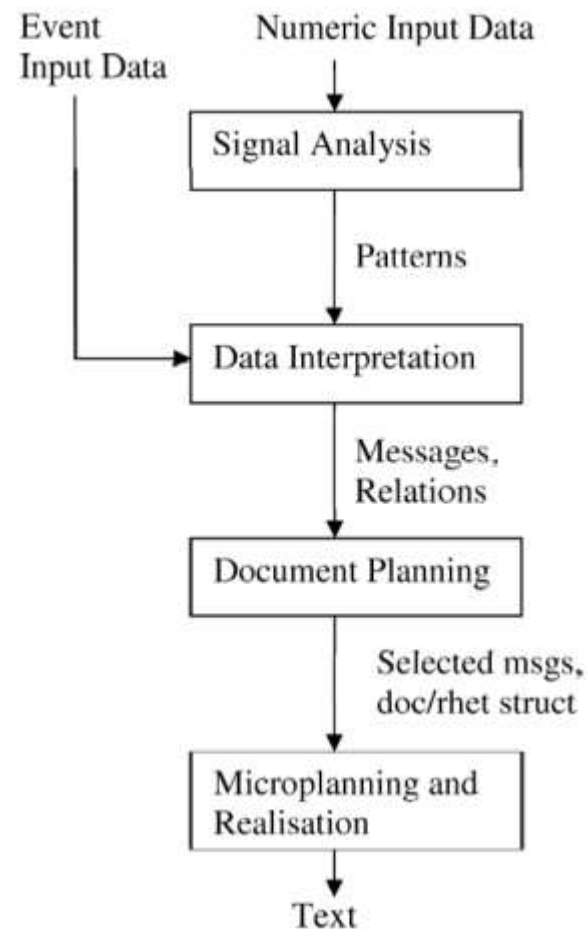
- 输入为基本模式与事件，通过对基本模式和输入事件进行分析，推断出更加复杂和抽象的消息，同时推断出它们之间的关系，最后输出高层消息以及消息之间的关系
- 例针对股票数据，如果跌幅超过某个值则可以创建一条消息，还需要检测消息之间的关系，例如因果关系、时序关系等



数据到文本生成系统框架

■ 文档规划模块(Document Planning)

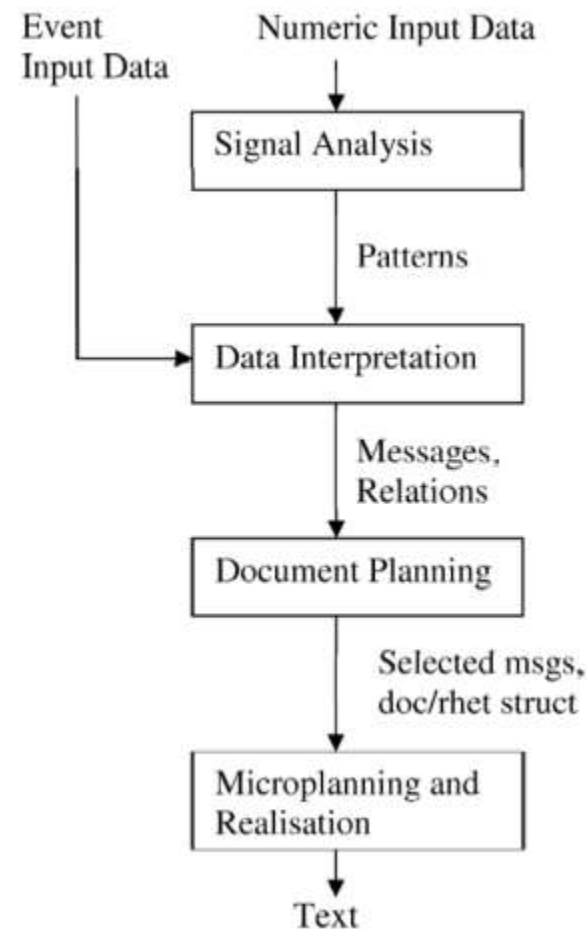
- 输入为消息及关系，分析决定哪些消息和关系需要在文本中提及，同时要确定文本的结构，最后输出需要提及的消息以及文档结构
- 信号分析与数据阐释模块会产生大量的消息、模式和事件，但文本通常长度受限，只能描述其中的一部分因此文档规划模块必须确定文本中需要说明的消息



数据到文本生成系统框架

■ 微规划与实现模块(Microplanning and Realisation)

- 输入为选中的消息及结构，通过自然语言生成技术输出最终的文本
- 主要涉及到对句子进行规划以及句子实现，要求最终实现的句子具有正确的语法、形态和拼写，同时采用准确的指代表达



数据到文本生成的应用

文本生成系统的应用领域：

- 天气预报领域的文本生成系统
- 针对空气质量的文本生成系统
- 针对财经数据的文本生成系统
- 面向医疗诊断数据的文本生成系统
- 基于体育数据生成文本摘要

数据到文本生成的应用

- 天气预报领域的文本生成技术应用最为成功
- FoG** 系统 能够从用户操作过的数据中生成双语天气预报文本
- SumTime** 系统能够生成海洋天气预报文本
- 英国阿伯丁大学的 **Anja Belz** 提出概率生成模型进行天气语言文本的生成
- Anja Belz** 和 **Eric Kow**进一步基于天气预报数据分析对比了多种数据到文本的生成系统

2. FORECAST 6 - 24 GMT, Wed 12-Jun 2002

WIND(KTS)

10M: W 8-13 backing SW by mid afternoon and S 10-15 by midnight.

50M: W 10-15 backing SW by mid afternoon and S 13-18 by midnight.

WAVES(M)

SIG HT: 0.5-1.0 mainly SW swell.

MAX HT: 1.0-1.5 mainly SW swell falling 1.0 or less mainly SSW swell by afternoon, then rising 1.0-1.5 by midnight.

PER(SEC)

WAVE PERIOD: Wind wave 2-4 mainly 6 second SW swell.

WINDWAVE PERIOD: 2-4.

SWELL PERIOD: 5-7.

WEATHER: Mainly cloudy with light rain showers becoming overcast around midnight.

VIS(NM): Greater than 10.

AIR TEMP(C): 8-10 rising 9-11 around midnight.

CLOUD(OKTAS/FT): 4-6 ST/SC 400-600 lifting 6-8 ST/SC 700-900 around midnight.

Figure 2. Forecast Text Produced by **SUMTIME**MOUSAM for the AM of 12-Jun 2002. The Wind part of the forecast has been generated from the data, shown in Table 1

数据到文本生成的应用

针对空气质量的文本生成系统

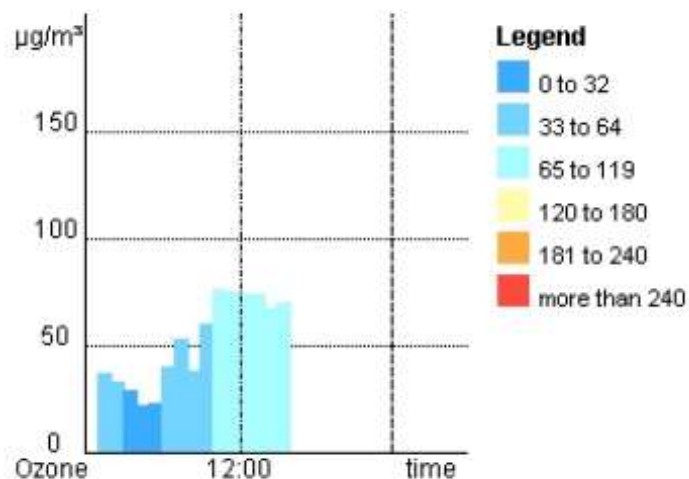


Figure 6: A sample English text

The ozone concentration ($70 \mu\text{g}/\text{m}^3$) is relatively low. As a result, no harmful effects to human health are expected. Between 4 AM and 10 AM, the ozone concentration increased considerably from 22 to 76. The current ozone concentration ($70 \mu\text{g}/\text{m}^3$) is close to the highest of $76 \mu\text{g}/\text{m}^3$ (at 10 AM). The lowest was $22 \mu\text{g}/\text{m}^3$ (at 4 AM).

数据到文本生成的应用

针对财经数据的文本生成系统

- 腾讯Dreamwriter，对数据进行学习，生成写作手法，进行新闻报道写作
- Narrative Science公司的生成系统可以自动生成新闻



腾讯Dreamwriter写稿机器人



Dreamwriter写稿流程

数据到文本生成的应

面向医疗诊断数据的文本生成系统

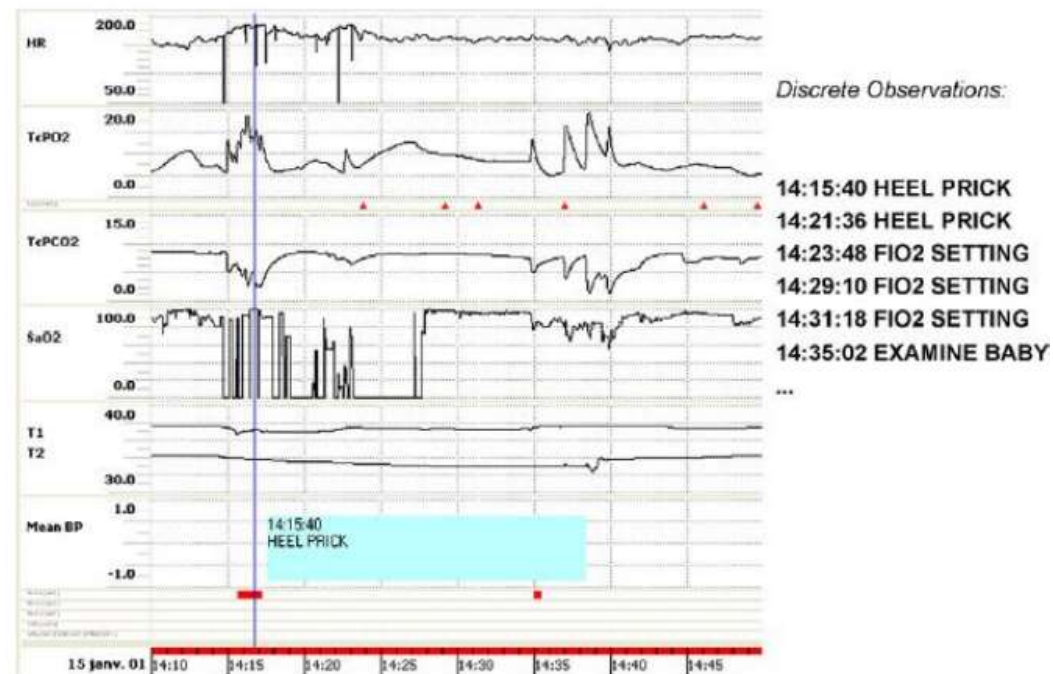


图 4.2: NICU 数据样例, 从上到下分别表示 HR, TcPO2, TcPCO2, SaO2, T1 & T2, and Mean BP

You saw the baby between 14:10 and 14:50. Heart Rate (HR) = 159. Core Temperature (T1) = 37.7. Peripheral Temperature (T2) = 34.3. Transcutaneous Oxygen (TcPO2) = 5.8. Transcutaneous CO2 (TcPCO2) = 8.5. Oxygen Saturation (SaO2) = 89.

Over the next 30 minutes T1 gradually increased to 37.3.

By 14:27 there had been 2 successive desaturations down to 56. As a result, Fraction of Inspired Oxygen (FIO2) was set to 45%. Over the next 20 minutes T2 decreased to 32.9. A heel prick was taken. Previously the spo2 sensor had been re-sited.

At 14:31 FIO2 was lowered to 25%. Previously TcPO2 had decreased to 8.4. Over the next 20 minutes HR decreased to 153.

By 14:40 there had been 2 successive desaturations down to 68. Previously FIO2 had been raised to 32%. TcPO2 decreased to 5.0. T2 had suddenly increased to 33.9. Previously the spo2 sensor had been re-sited. The temperature sensor was re-sited.

图 4.3: BT-45 系统生成的对应文本 [Portet et al., 2009]^[87]

数据到文本生成的应用

工业界成立了多家从事文本生成的公司，能够为多个行业基于行业数据生成行业报告或新闻报道等

ARRIA

narrative  science

ai AUTOMATED
INSIGHTS®

INVESTING 7/07/2015 1:00下午 | 332 views

Earnings for Alcoa Projected to Rise

By Narrative Science

[+ Comment Now](#) [+ Follow Comments](#)

Wall Street is high on **Alcoa**, expecting it to report earnings that are up 28% from a year ago when it reports its second-quarter earnings on Wednesday, July 8, 2015. The consensus estimate is 23 cents per share, up from earnings of 18 cents per share a year ago.

The consensus estimate has fallen over the past three months, from 27 cents. Analysts are expecting earnings of 95 cents per share for the fiscal year. Analysts look for revenue to decrease 1% year-over-year to \$5.79 billion for the quarter, after being \$5.84 billion a year ago. For the year, revenue is projected to roll in at \$23.63 billion.

Revenue dropped year-over-year in the first quarter, ending a two-quarter streak of growing revenue.

Alcoa is a global producer of aluminum. It is mainly engaged in the production and management of primary aluminum, fabricated aluminum, and alumina combined. It is actively involved in a range of industries, including technology, mining, smelting, and recycling. Kaiser Aluminum Corp., also in the metal mining industry, will report earnings on Wednesday, July 22, 2015. Analysts are expecting earnings of \$1.19 per share for Kaiser Aluminum, up 13% from last year's earnings of \$1.05 per share. Other companies in the metal mining industry with upcoming earnings release dates include: Noranda Aluminum Holding and Aluminum Corp. of China Limited.

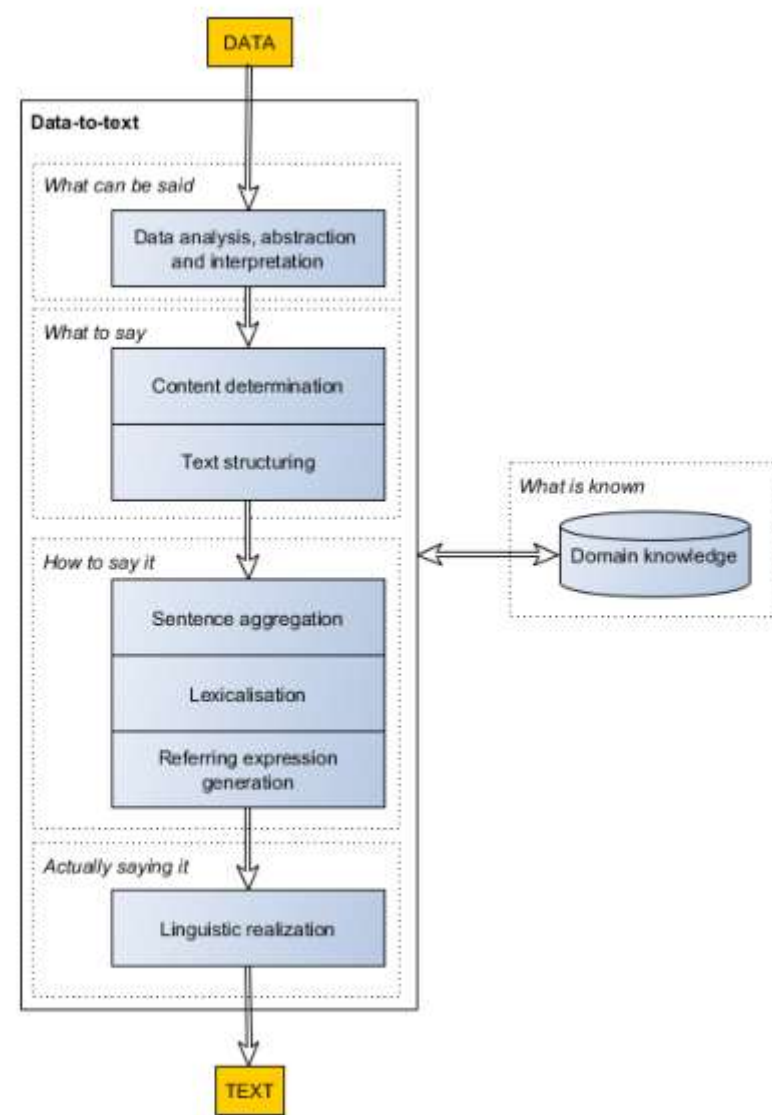
Earnings estimates provided by Zacks.

图 4.4: NarrativeScience 自动生成的样例新闻

数据到文本生成的研究

数据到文本生成的医疗领域

- report automation
- clinical decision support
- behaviour change
- patient engagement
- patient assistance



数据到文本生成的研究前沿

数据到文本的新闻领域

- Munezero提出一种数据驱动的自动化新闻系统

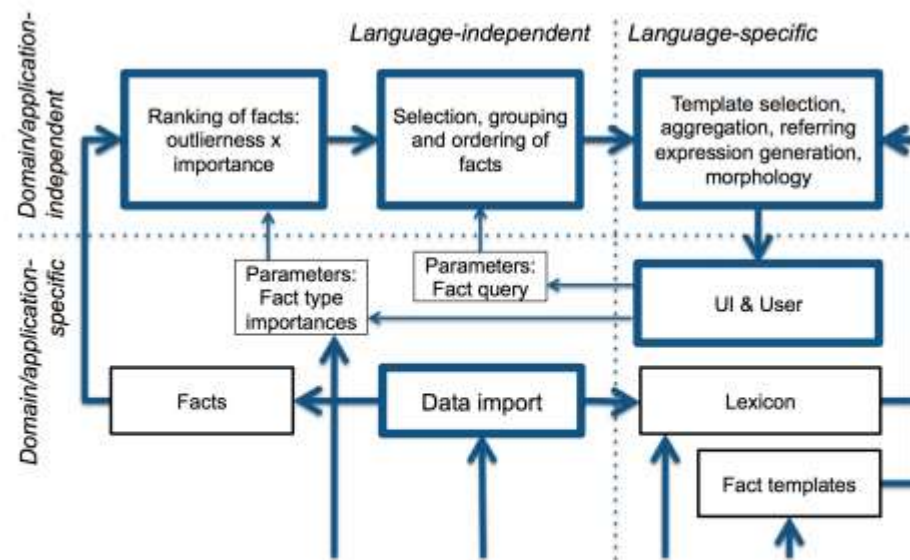


Figure 1: Overview of the architecture. Thick boxes represent software components and thin boxes data structures.

小结

数据到文本生成的定义

数据到文本生成系统框架

- 三阶段流水

数据到文本生成的应用

- 天气预报、财经、医疗等

数据到文本生成的研究前沿

- 医疗、新闻

三、文本到文本的生成

文本到文本的生成定义

对联自动生成

诗歌自动生成

总结与展望

定义

定义：对给定文本进行变换和处理从而获得新文本的技术

应用：

- 对联自动生成
- 诗歌自动生成
- 作文自动生成
- 对话生成
- 等等

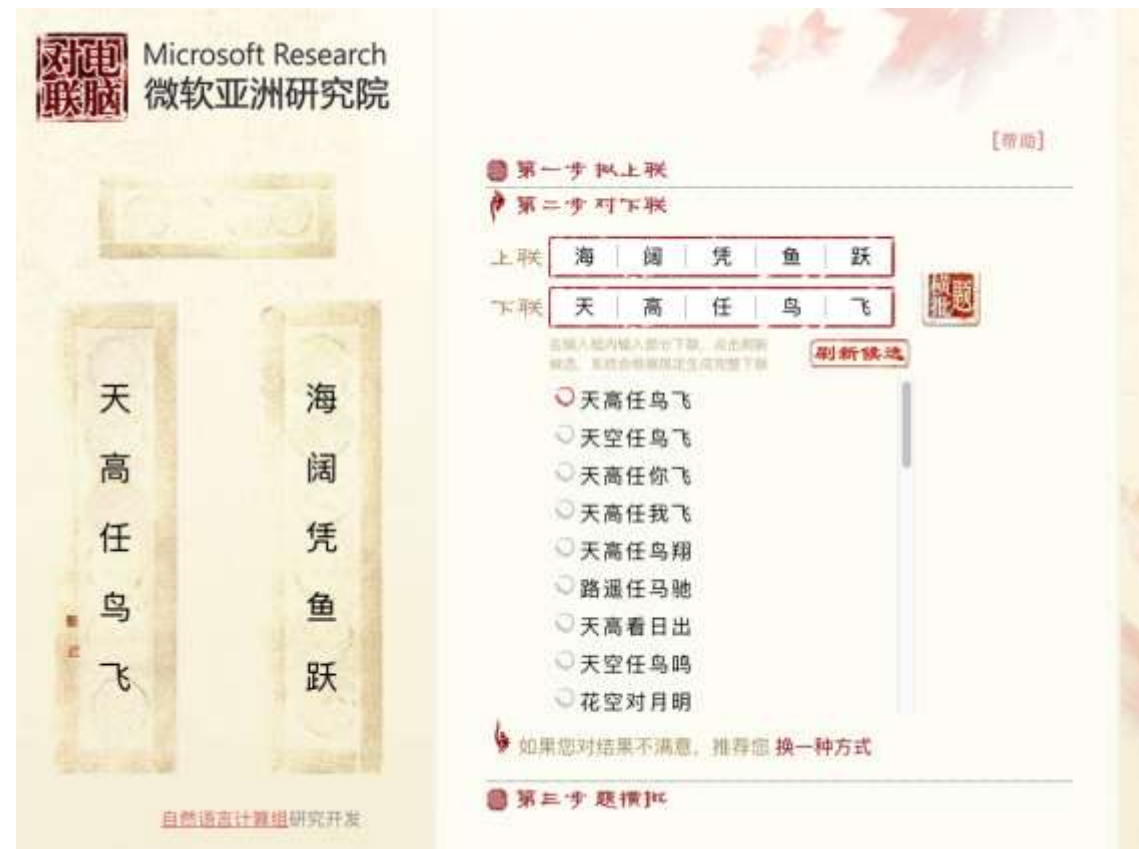
对联自动生成

对联特点：上联(FS, first sentence)和下联(SS, second sentence)

对联自动生成：/给出上联，自动生成下联

海	阔	凭	鱼	跃
sea	wide	allow	fish	jump
天	高	任	鸟	飞
sky	high	permit	bird	fly

微软对联：<http://duilian.msra.cn>



对联特点

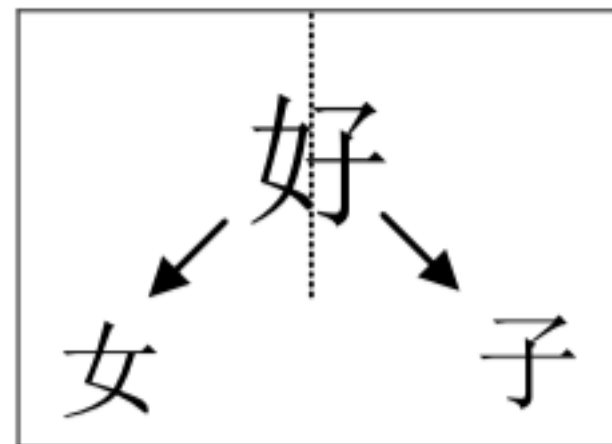
1. 上下联具有相同的长度，对应位置的分词一致
2. 上下联的音调协调。**FS**的最后一个字是仄音，**SS**的最后一个字是平音
3. 上下联对应的词具有相同的词性和特点
4. 上下联内容相关，但不重复

海	阔	凭	鱼	跃
sea	wide	allow	fish	jump
天	高	任	鸟	飞
sky	high	permit	bird	fly

对联特点

5. 上下联具有相同的写作风格，如 **FS** 中有重复的字、词或发音，则 **SS** 中对应位置具有相同的重复；**FS** 中具有字的分解，**SS** 中对应位置也应该有对应字的分解

有	女	有	子	方	称	好
have	daughter	have	son	so	call	good
缺	鱼	缺	羊	敢	叫	鲜
lack	fish	lack	mutton	dare	call	delicious



对联自动生成

对联自动生成步骤：

1. 基于 “_____” 的对联自动生成{请填空}

“天对地,雨对风。大陆对长空。山花对海树,赤日对苍穹”

1. **Phrase-based SMT Model** 生成出**N-Best** 候选
2. 基于语言学的筛选
3. 基于其他特征的重排序

对联自动生成模型

Phrase-based SMT Model

- 上联表示 $F = \{f_1, f_2, \dots, f_n\}$, 下联表示 $S = \{s_1, s_2, \dots, s_n\}$, 其中 f_i, s_i 是对联中对应第 i 个汉字
- 生成对联, 目标最大化 $p(S/F)$
- 采用 phrase-based log-linear model

$$\begin{aligned} S^* &= \arg \max_S p(S | F) \\ &= \arg \max_S \sum_{i=1}^M \lambda_i \log h_i(S, F) \end{aligned}$$

其中 $h_i(S, F)$ 为特征函数, M 为特征函数的个数, λ_i 是估计量

对联自动生成模型

应用于Phrase-based SMT Model的特征：

S, F 切分为短语后分别表示为 $\bar{s}_1, \dots, \bar{s}_l$ 和 $\bar{f}_1, \dots, \bar{f}_l$

1. Phrase translation model
2. Inverted phrase translation model
3. Lexical weight
4. Inverted lexical weight
5. 下联的语言模型得分(character-based trigram)

$h_1(S, F) = \prod_{i=1}^l p(\bar{f}_i \bar{s}_i)$	Phrase translation model
$h_2(S, F) = \prod_{i=1}^l p(\bar{s}_i \bar{f}_i)$	Inverted phrase translation model
$h_3(S, F) = \prod_{i=1}^l p_w(\bar{f}_i \bar{s}_i)$	Lexical weight
$h_4(S, F) = \prod_{i=1}^l p_w(\bar{s}_i \bar{f}_i)$	Inverted lexical weight
$h_5(S, F) = p(S)$	Language model

基于语言学的筛选

1. Repetition filter

检查上联和下联对应位置是否有重复的字(词)

2. Pronunciation repetition filter

检查上联和下联对应位置字的发音

3. Character decomposition filter

检查上联和下联字的分解是否对应

4. Phonetic harmony filter

根据对联下联发音特点，过滤最后一词发音不符的下联候选

基于其他特征重排序

FS: 海阔凭鱼跃

SS1: 天高任鸟飞

SS2: 天高任狗叫

SS1和SS2具有相同的得分(SMT、语言模型), 但”天高”和“狗叫”放在一起没有意义

1. Mutual information (MI) score

$$MI(S) = \sum_{i=1}^{n-1} \sum_{j=i+1}^n I(s_i, s_j) = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \log \frac{p(s_i, s_j)}{p(s_i) p(s_j)}$$

基于其他特征重排序

海阔凭鱼跃，天高任鸟飞

FS: “海”和“阔”，“海”和“鱼”，“鱼”和“跃”具有较强联系

SS: “天”和“高”，“天”和“鸟”，“鸟”和“飞”具有较强联系

2. MI-based structural similarity

$$F = \{f_1, f_2, \dots, f_n\}$$

$$V_f = \{v_{12}, v_{13}, \dots, v_{1n}, v_{23}, \dots, v_{n-1n}\} \quad (V_{ij} \text{ 是 } f_i, f_j \text{ 的互信息得分})$$

$$MISS(F, S) = \cos(V_f, V_s) = \frac{V_f \bullet V_s}{|V_f| \times |V_s|}$$

Reranking model: Ranking SVM

生成对联的评价方法

指标：BLEU, human evaluation

1. Bleu与人工评价具有0.92相关性
2. 所选特征对于生成对联的有效性
3. 人工评价自动生成对联中可接受的比例

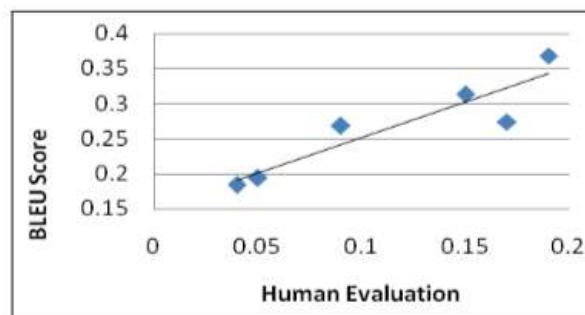


Figure 4: BLEU Predicts Human Judgments.

	Features	BLEU
Phrase-based SMT Model	Phrase TM(PTM) + LM	0.276
	+ Inverted PTM	0.282
	+ Lexical Weight (LW)	0.315
	+ Inverted LW	0.348
Ranking SVM	+ Mutual information (MI)	0.356
	+ MI-based structural similarity	0.361

Table 3. Feature Evaluation.

	Top-1	Top-10
Top-n inclusion rate	0.21	0.73

Table 4. Overall Performance Evaluation.

诗句自动生成

输入代表诗主题的词, 逐句生成直到末句, 最后诗题



Microsoft
Research
微软亚洲研究院

作诗一首 (重新作诗)

白日依山尽

黄河入海流

欲穷千里目

更上一层楼

寫詩題

推荐诗句

草稿箱(0)

主题词: + 确定

五言诗句

更上一层楼

难空一边隙

不富一中垢

难风一中奏

难富一边游

无限百中秋

将老一中秋

不富万年油

无空百中游

不限百中首

将富百中舟

不富一中愁

无限几间幽

思空百年厚

无限万中菱

思风四年头

上一页 第 1 页v 下一页

诗句特点

1. 严格的声调模式(平仄)

五言绝句和七言绝句具有4中常见的平仄结构，右侧为其中一种

* + - - +
- - + + -
* - - + +
* + + - -

+代表仄，-代表平，*代表任意一种

2. 押韵

押韵的字具有具有相同元音结尾，诗句第二句和第四句的结尾押韵

3. 结构限制

绝句具有“起, 承, 转, 合”的结构

登鹳雀楼
On the Stork Tower
王之涣

白日依山尽, (- + - - +)
<i>white sunlight along hill fade</i>
黄河入海流。(- - + + -)
<i>Yellow River into sea flow</i>
欲穷千里目, (+ - - + +)
<i>wish exhaust thousand mile eyesight</i>
更上一层楼。(+ + + - -)
<i>More up one story tower</i>

诗句自动生成

1. 基于模板生成第一句

- a) 用户从已经聚类的短语对应的词中选择关键词，根据关键词得到候选短语
- b) 所有候选短语放入可以满足平仄模式的所有可能位置，构成短语格子
- c) 采用语言模型计算格子中所有路径的得分
- d) **Forward-Viterbi-Backward-A*** 算法选择**N-best**候选句子

2. **SMT model**生成四句诗

- a) 上句生成下句作为翻译任务，采用**SMT model**(对联自动生成的**SMT**)
- b) 三个不同的**SMT**系统，分别生成第二句、第三句和第四句
- c) **Coherence model**将之前的句子中的关键词用于**SMT** 系统，排序生成句子候选

诗句自动生成

Coherence Model

$$\begin{aligned} S^* &= \arg \max_S p(S | F) \\ &= \arg \max_S \sum_{i=1}^M \lambda_i \log h_i(S, F) \end{aligned}$$

1. 上述用于**SMT**的公式只考虑到相邻绝句的信息，而诗四个句子都要连贯
2. **Coherence Model** 计算下一句与之前生成所有的句子的互信息，当做**SMT**系统的第六个特征

$$h_6(S, F) = \sum_{i,j} MI(s_i, s_j) = \sum_{i,j} \log \frac{p(s_i, s_j)}{p(s_i)p(s_j)}$$

s_i, s_j 代表SetA和SetB中的字，SetA由已生成的诗句的字组成，SetB由下一句的字组成

$h_1(S, F) = \prod_{i=1}^I p(\bar{f}_i \bar{s}_i)$	Phrase translation model
$h_2(S, F) = \prod_{i=1}^I p(\bar{s}_i \bar{f}_i)$	Inverted phrase translation model
$h_3(S, F) = \prod_{i=1}^I p_w(\bar{f}_i \bar{s}_i)$	Lexical weight
$h_4(S, F) = \prod_{i=1}^I p_w(\bar{s}_i \bar{f}_i)$	Inverted lexical weight
$h_5(S, F) = p(S)$	Language model

Table 1. Features Used in the SMT Model.

展望

展望：

- 自动生成对话
- 自动生成作文
-

补充学习：图像到文本的生成

定义

三阶段的流水线模式

相关研究现状

数据集与评价指标

小结

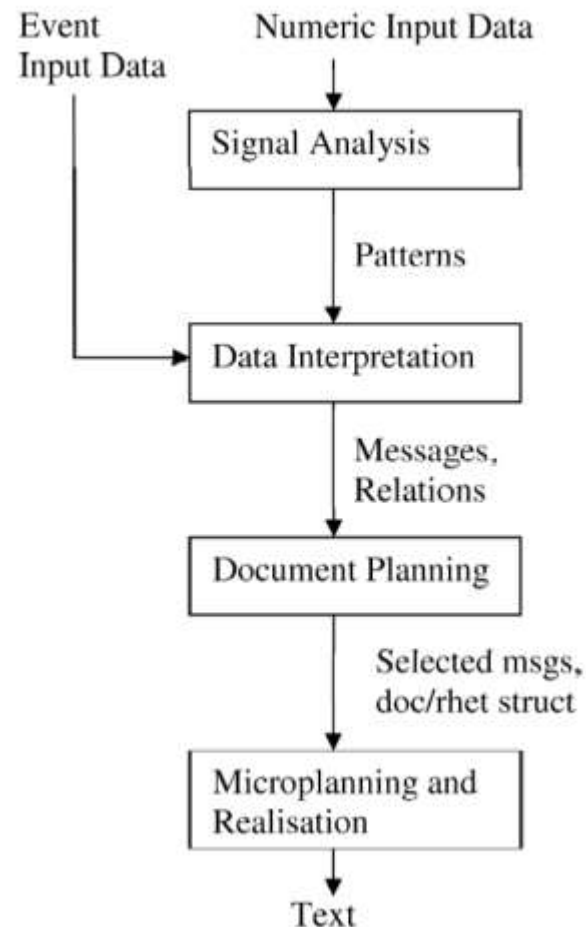
定义及类别

- ◆ 根据给定的图像生成描述该图像内容的自然语言文本
 - ◆ 新闻图像附带的标题
 - ◆ 医学图像附属的说明
 - ◆ 儿童教育中常见的看图说话
 - ◆ 以及用户在微博等互联网应用中上传图片时提供的说明文字
- ◆ 根据文本的详细程度及长度
 - ◆ 分为图像标题自动生成
 - ◆ 图像说明自动生成

图像到文本的生成技术

根据图像内容理解的特点，包括三阶段

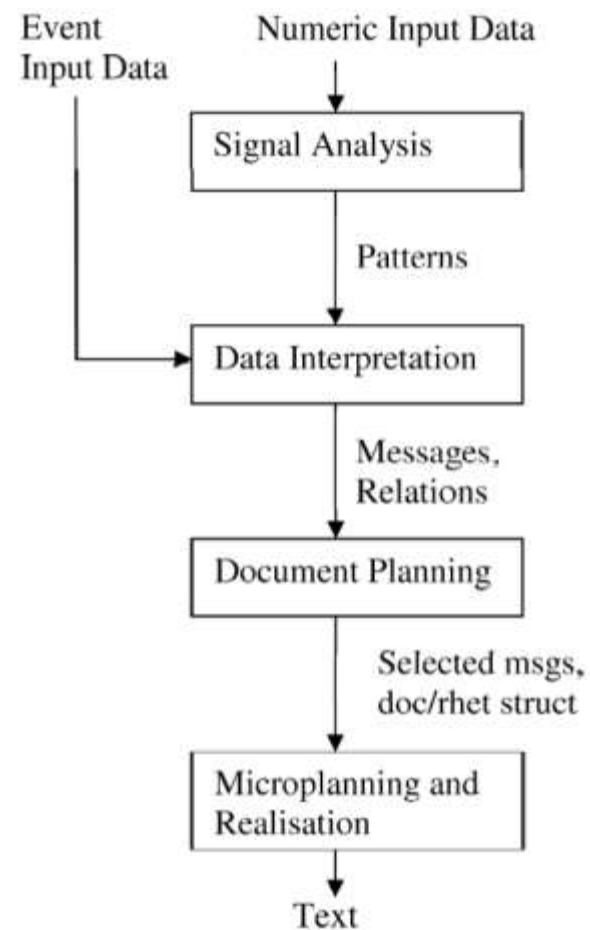
- 内容抽取
- 句子内容选择
- 句子实现



三阶段的流水线模式

内容抽取方面，

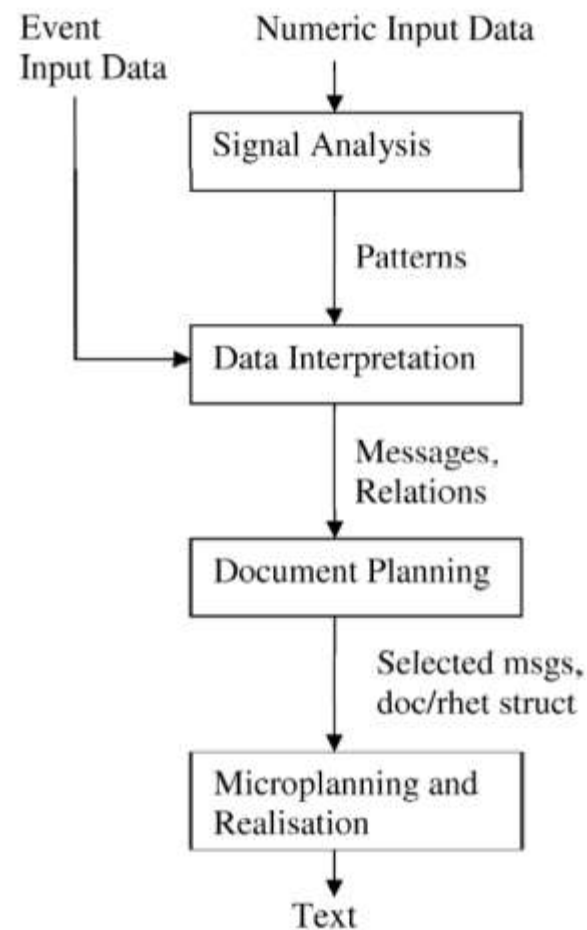
- 需要从图像中抽取物体、方位、动作、场景等概念，其中物体可以具体定位到图像中的某一具体区域，而其他概念则需要语义标引。这部分主要依靠模式识别和计算机视觉技术



三阶段的流水线模式

句子内容选择方面

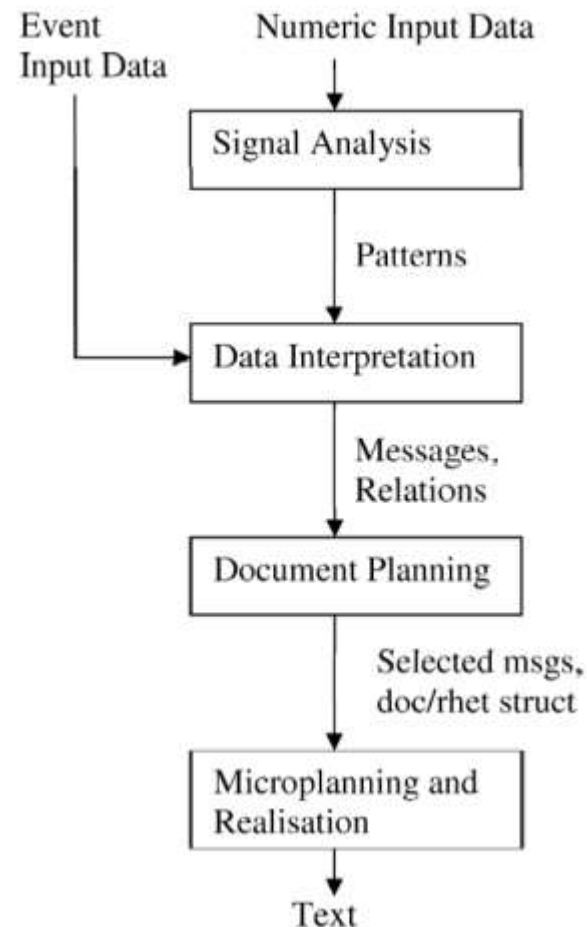
- 需要依据应用场景，选择最重要（如图像画面中最突出的，或与应用场景最相关的），且意义表述连贯的概念。这部分需要综合运用计算机视觉与自然语言处理技术



三阶段的流水线模式

句子实现部分

- 根据实际应用特点选取适当的表述方式将所选择的概念梳理为合乎语法习惯的自然语言句子。这部分主要依靠自然语言处理技术



相关研究现状

依靠三阶段流水线模式

- 基于图像描述模板
- 基于概率图模型

依靠计算机视觉提取图像中物体及定位技术

- 基于概率图和语言模型
- 基于核函数的典型关联分析

图像语义标注与自然语言句子生成联合建模

- 基于多模态**m-RNN**
- 基于视觉中注意力机制

基于图像描述模板

图像被细致的分割并标注为物体及其组成部分，以及图像所表现的场景，并在此基础上选择与场景相关的描述模板，将物体识别的结果填充入模板得到图像的描述文字

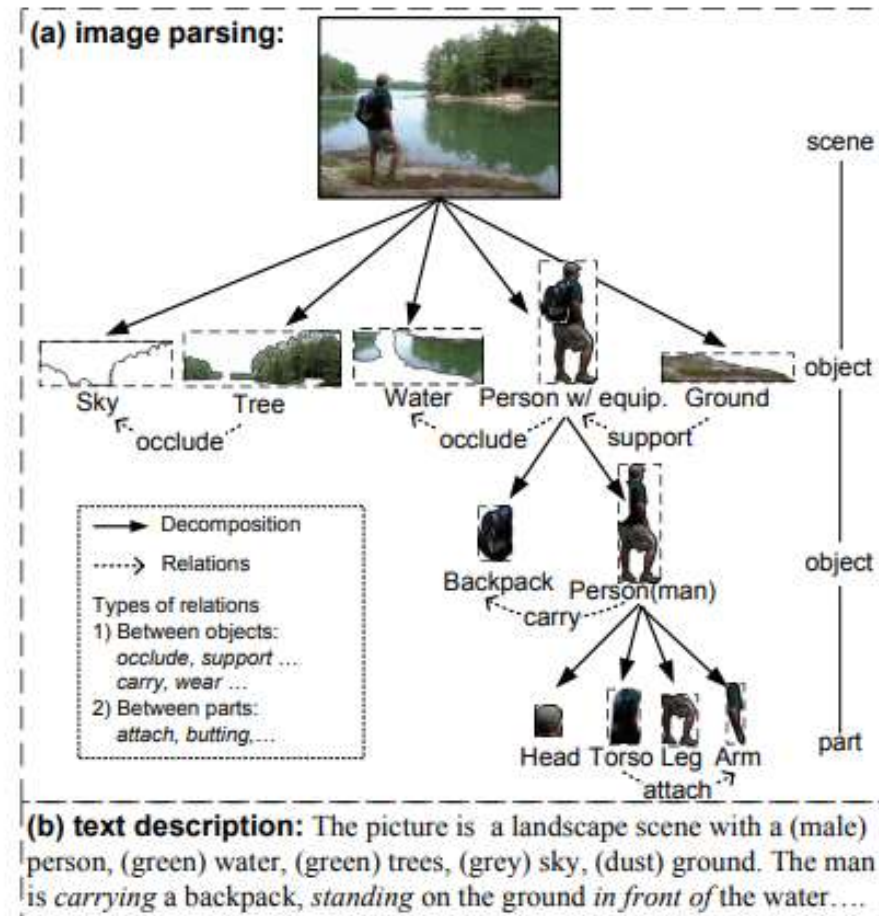


Fig. 1. Two major tasks of the I2T framework: (a) image parsing and (b) text description. See text for more details.

基于图像描述模板

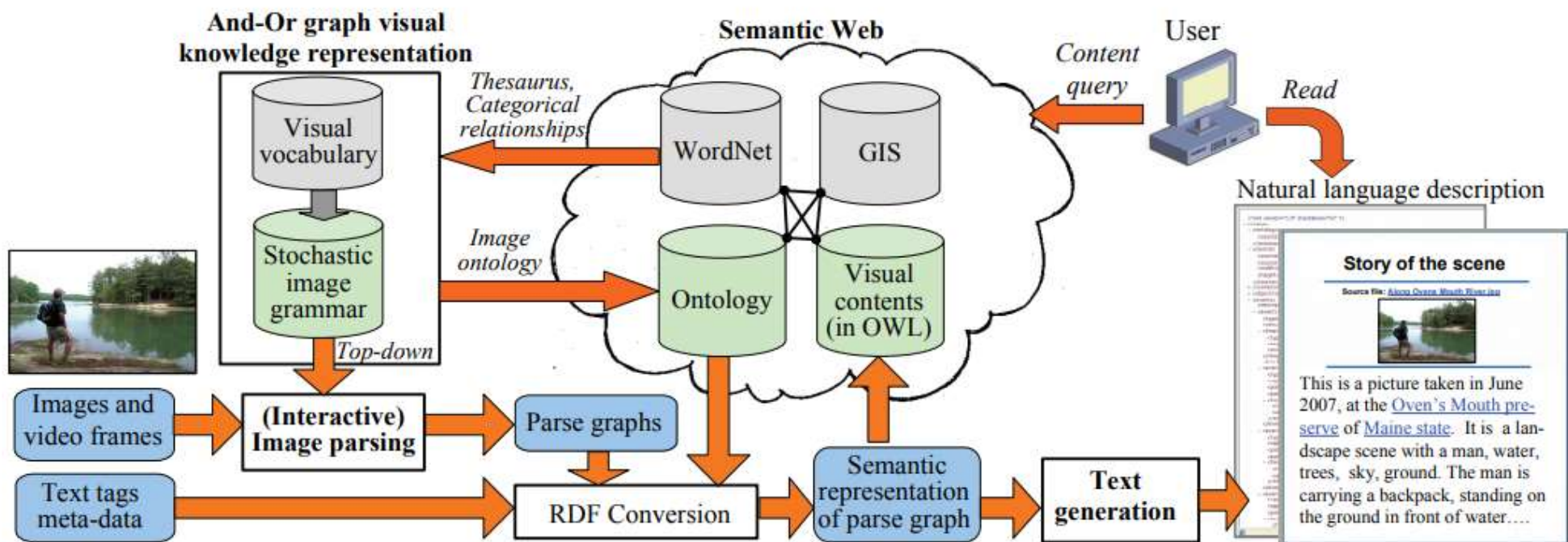


Fig. 2. Diagram of the I2T framework. Four key components (as highlighted by bold fonts) are: (1) An image parsing engine that converts input images or video frames into parse graphs. (2) An And-or Graph visual knowledge representation that provides top-down hypotheses during image parsing and serves as an ontology when converting parse graphs into semantic representations in RDF format. (3) A general knowledge base embedded in the Semantic Web that enriches the semantic representations by interconnecting several domain specific ontologies. (4) A text generation engine that converts semantic representations into human readable and query-able natural language descriptions.

基于概率图模型

采用概率图模型对文本信息和图像信息同时建模，并从新闻图片所在的文字报道中挑选合适的关键词作为体现图像内容的关键词，并进而利用语言模型规则链接为基本合乎语法





<p>Thousands of Tongans have attended the funeral of King Taufa'ahau Tupou IV, who died last week at the age of 88. Representatives from 30 foreign countries watched as the king's coffin was carried by 1,000 men to the official royal burial ground.</p>		<p>King Tupou, who was 88, died a week ago.</p>
<p>Contaminated Cadbury's chocolate was the most likely cause of an outbreak of salmonella poisoning, the Health Protection Agency has said. About 36 out of a total of 56 cases of the illness reported between March and July could be linked to the product.</p>		<p>Cadbury will increase its contamination testing levels.</p>
<p>A Nasa satellite has documented startling changes in Arctic sea ice cover between 2004 and 2005. The extent of "perennial" ice declined by 14%, losing an area the size of Pakistan or Turkey. The last few decades have seen ice cover shrink by about 0.7% per year.</p>		<p>Satellite instruments can distinguish "old" Arctic ice from "new".</p>
<p>A third of children in the UK use blogs and social network websites but two thirds of parents do not even know what they are, a survey suggests. The children's charity NCH said there was "an alarming gap" in technological knowledge between generations.</p>		<p>Children were found to be far more internet-wise than parents.</p>

Table 1: Each entry in the BBC News database contains a document an image, and its caption.

依靠计算机视觉提取图像中物体及定

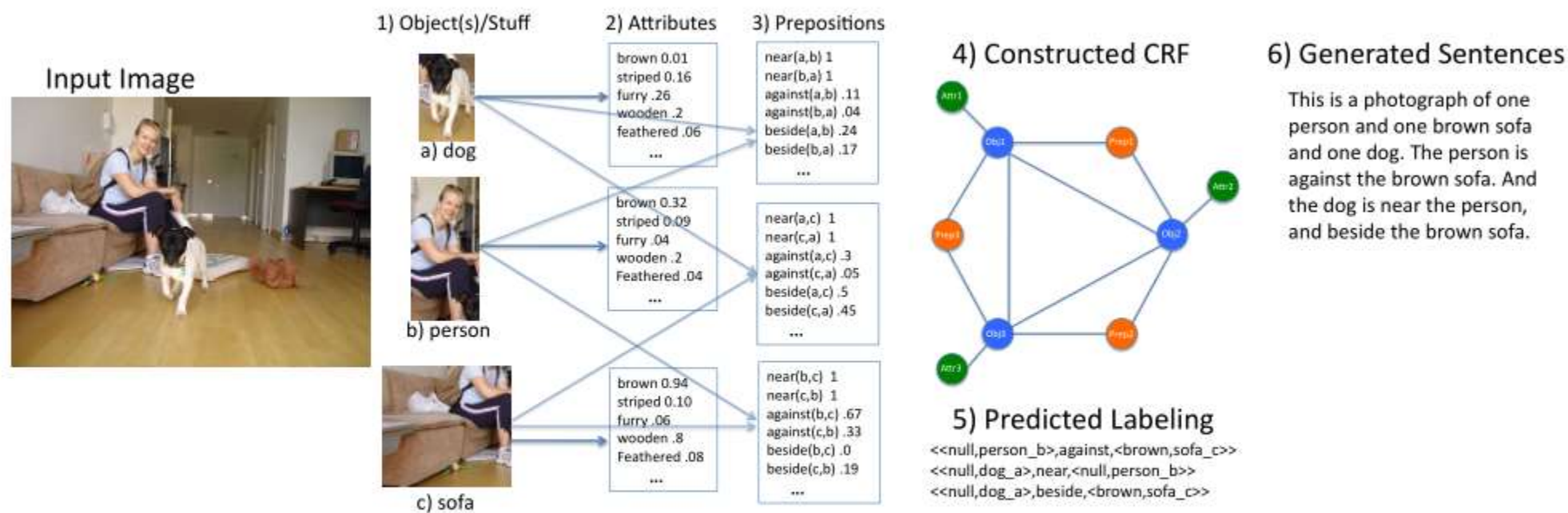


Figure 2. System flow for an example image: 1) object and stuff detectors find candidate objects, 2) each candidate region is processed by a set of attribute classifiers, 3) each pair of candidate regions is processed by prepositional relationship functions, 4) A CRF is constructed that incorporates the unary image potentials computed by 1-3, and higher order text based potentials computed from large document corpora, 5) A labeling of the graph is predicted, 6) Sentences are generated based on the labeling.

计算机视觉提取图像中物体及定位

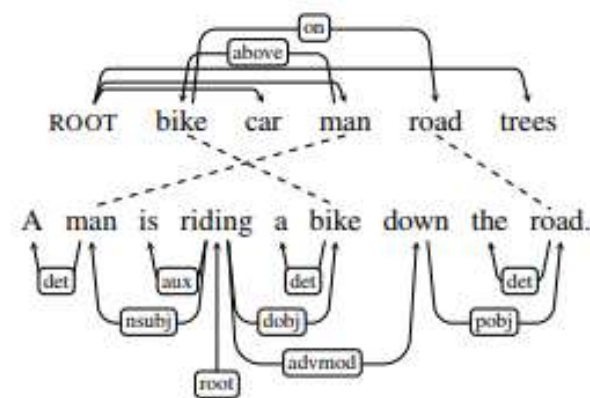
依靠计算机视觉领域现有的物体识别技术从图像中抽取物体（包括人物、动物、花草、车、桌子等常见的物体类型），并对其定位以获得物体之间的上下位关系，进而依赖概率图模型和语言模型选取适当的描述顺序将这些物体概念、介词短语块串联成完整的句子



(a)

A man is riding a bike down the road.
A car and trees are in the background.

(b)



(c)

Figure 1: (a) Image with regions marked up: BIKE, CAR, MAN, ROAD, TREES; (b) human-generated image description; (c) visual dependency representation expressing the relationships between MAN, BIKE, and ROAD aligned to the syntactic dependency parse of the first sentence in the human-generated description (b).

基于核函数的典型关联分析

使用**Kernel Canonical Correlation Analysis(KCCA)**将图片和文本映射到共享的潜在语义空间



A man is doing tricks on a bicycle on ramps in front of a crowd.
A man on a bike executes a jump as part of a competition while the crowd watch
A man rides a yellow bike over a ramp while others watch.
Bike rider jumping obstacles.
Bmx biker jumps off of ramp.

Figure 1: An example of an image from the Flickr 8K dataset. Each of the captions literally describe what is being depicted in the photograph while also mentioning different entities and exhibiting linguistic variation

图像语义标注与自然语言句子生成联合建模

在图像端采用多层深度卷积神经网络（**Deep Convolution Neural Network, DCNN**）对图像中的物体概念进行建模

在文本端采用循环神经网络（**Recurrent Neural Network, RNN**）或递归神经网络（**Recursive Neural Network**）对自然语言句子的生成过程进行建模

多模态M-RNN模型

Mao 等人通过 DCNN 得到的图像信息与文本信息融合到同一个循环神经网络 (m-RNN) 中，将图像信息融入到了自然语言句子生成的序列过程

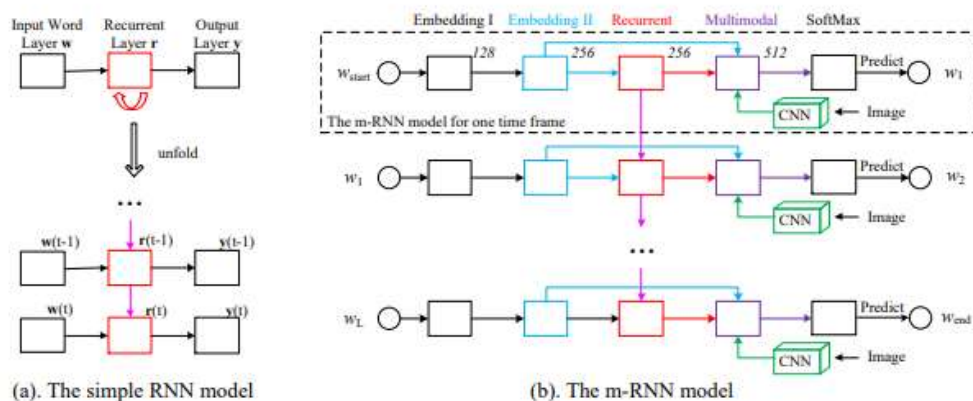


Figure 2: Illustration of the simple Recurrent Neural Network (RNN) and our multimodal Recurrent Neural Network (m-RNN) architecture. (a). The simple RNN. (b). Our m-RNN model. The inputs of our model are an image and its corresponding sentence descriptions. w_1, w_2, \dots, w_L represents the words in a sentence. We add a start sign w_{start} and an end sign w_{end} to all the training sentences. The model estimates the probability distribution of the next word given previous words and the image. It consists of five layers (i.e. two word embedding layers, a recurrent layer, a multimodal layer and a softmax layer) and a deep CNN in each time frame. The number above each layer indicates the dimension of the layer. The weights are shared among all the time frames. (Best viewed in color)

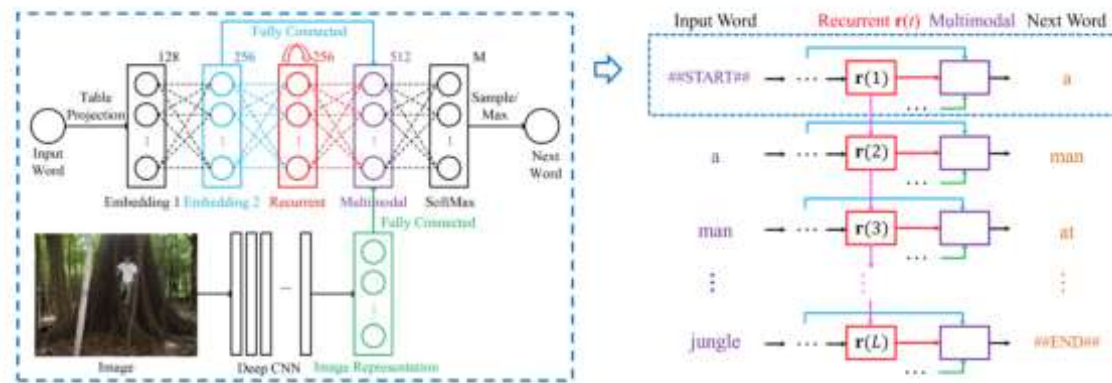
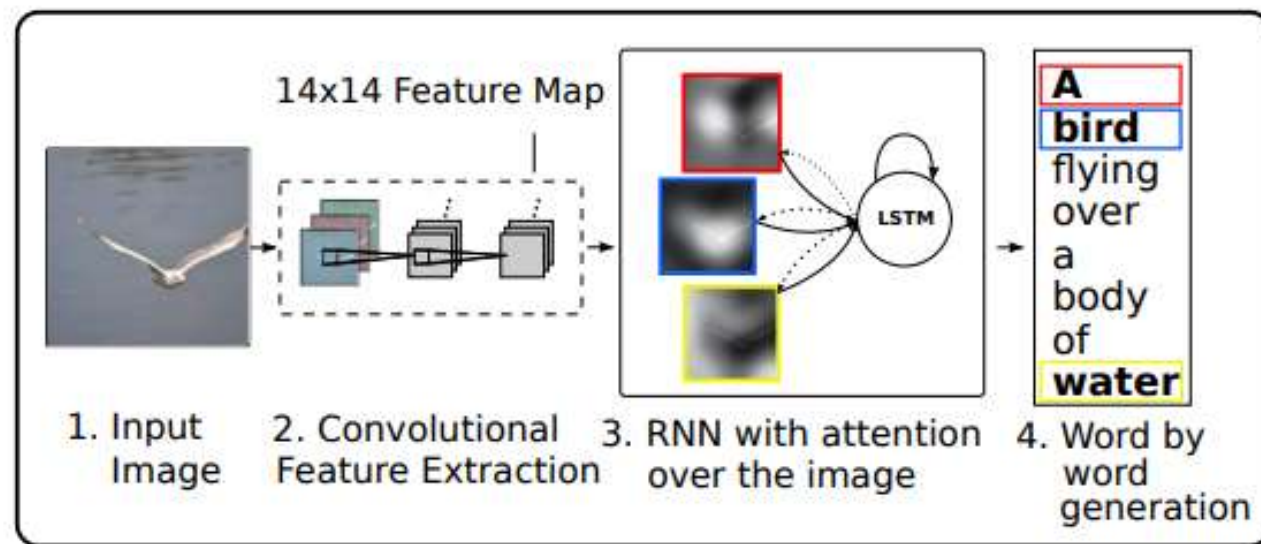


图 5.2: 多模态 m-RNN 模型^[100]

视觉“注意”引导的图像标题生成

Xu等人提出利用计算机视觉领域中的“注意”(Attention)机制来促进词语和图像块之间的对齐，从而在句子生成过程中，模拟人视觉的“注意”转移过程能够与词语序列的生成过程相互促进，使生成的句子更符合人的表述习惯

Figure 1. Our model learns a words/image alignment. The visualized attentional maps (3) are explained in Sections 3.1 & 5.4



基于卷积神经网络 CNN 和多示例学习

微软的研究人员利用卷积神经网络 CNN 和多示例学习 (Multiple Instance Learning, MIL) 对图像建模, 并利用判别式语言模型生成候选句子, 并采用统计机器翻译研究中经典的最小误差率训练 (Minimum Error Rate Training, MERT) 来发掘文本和图像层面的特征对候选句子进行排序



D-ME+DMSM

a plate with a sandwich and a cup of coffee

MRNN

a close up of a plate of food

D-ME+DMSM+MRNN

a plate of food and a cup of coffee

k-NN

a cup of coffee on a plate with a spoon



D-ME+DMSM

a black bear walking across a lush green forest

MRNN

a couple of bears walking across a dirt road

D-ME+DMSM+MRNN

a black bear walking through a wooded area

k-NN

a black bear that is walking in the woods



D-ME+DMSM

a gray and white cat sitting on top of it

MRNN

a cat sitting in front of a mirror

D-ME+DMSM+MRNN

a close up of a cat looking at the camera

k-NN

a cat sitting on top of a wooden table

Table 2: Example generated captions.

数据集

	Images	Texts	Judgments	Objects
Pascal1K (Rashtchian et al., 2010)	1,000	5	No	Partial
VLT2K (Elliott & Keller, 2013)	2,424	3	Partial	Partial
Flickr8K (Hodosh & Hockenmaier, 2013)	8,108	5	Yes	No
Flickr30K (Young et al., 2014)	31,783	5	No	No
Abstract Scenes (Zitnick & Parikh, 2013)	10,000	6	No	Complete
IAPR-TC12 (Grubinger et al., 2006)	20,000	1–5	No	Segmented
MS COCO (Lin et al., 2014)	164,062	5	Collected	Partial
BBC News (Feng & Lapata, 2008)	3,361	1	No	No
SBU1M Captions (Ordonez et al., 2011)	1,000,000	1	Collected ⁷	No
Déjà-Image Captions (Chen et al., 2015)	4,000,000	Varies	No	No

Table 2: Image datasets for the automatic description generation models. We have split the overview into image *description* datasets (top) and *caption* datasets (bottom) – see the main text for an explanation of this distinction.

数据集



1. One jet lands at an airport while another takes off next to it.
2. Two airplanes parked in an airport.
3. Two jets taxi past each other.
4. Two parked jet airplanes facing opposite directions.
5. two passenger planes on a grassy plain

(a) Pascal1K⁸



1. There are several people in chairs and a small child watching one of them play a trumpet
2. A man is playing a trumpet in front of a little boy.
3. People sitting on a sofa with a man playing an instrument for entertainment.

(b) VLT2K⁹

数据集



1. A man is snowboarding over a structure on a snowy hill.
2. A snowboarder jumps through the air on a snowy hill.
3. a snowboarder wearing green pants doing a trick on a high bench
4. Someone in yellow pants is on a ramp over the snow.
5. The man is performing a trick on a snowboard high in the air.

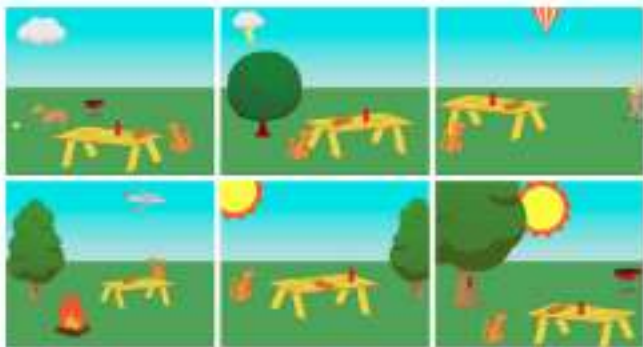
(c) Flickr8K¹⁰



1. a yellow building with white columns in the background
2. two palm trees in front of the house
3. cars are parking in front of the house
4. a woman and a child are walking over the square

(d) IAPR-TC12¹¹

数据集



1. A cat anxiously sits in the park and stares at an unattended hot dog that someone has left on a yellow bench

(e) Abstract Scenes¹²



1. A blue smart car parked in a parking lot.
2. Some vehicles on a very wet wide city street.
3. Several cars and a motorcycle are on a snow covered street.
4. Many vehicles drive down an icy street.
5. A small smart car driving in the city.

(f) MS COCO¹³

评价指标

- ✓ 与标准文本对比
 - ✓ 直接利用机器翻译自动评价指标
 - ✓ BLEU
 - ✓ Meteor
 - ✓ ROUGE

小结

定义

三阶段流水模型

- 内容抽取方面
- 句子内容选择方面
- 句子实现方面

相关研究现状

- 依靠三阶段流水线模式
- 依靠计算机视觉提取图像中物体及定位技术
- 图像语义标注与自然语言句子生成联合建模

数据集与评价指标