

机器翻译系统融合方法及其应用探究*

□ 何彦青 石崇德 于薇 张均胜 王惠临 / 中国科学技术信息研究所 北京 100038

摘要: 多机器翻译系统融合技术能够对不同机器翻译系统的输出结果有效地进行融合,产生更好的翻译性能,因此该技术成为机器翻译研究领域的一个热点问题。文章介绍了中国科学技术信息研究所(ISTIC)参加第七届全国机器翻译研讨会机器翻译评测的情况。本单位参加了英汉科技领域的机器翻译评测项目。文章阐述了本单位机器翻译系统的实现框架以及实施细节,并分析了它们在评测数据上的性能表现,最后对机器翻译系统融合方法目前的现状进行讨论,并对该系统融合方法进行总结和展望。

关键词: 机器翻译,自然语言处理,系统融合

DOI: 10.3772/j.issn.1673-2286.2011.12.005

1 引言

机器翻译经过半个多世纪的发展和研究,各种翻译方法不断涌现,呈现出百家争鸣的态势,基于规则的机器翻译和基于统计的机器翻译相继登场,各自占据一方天地。越来越多不同类型的翻译系统出现,它们的翻译模型差异使得翻译结果呈现出不同的翻译特征。即使是在同一个翻译系统中,为了平衡翻译效果和翻译效率之间的关系,其解码器在寻优过程中加入的多种剪枝策略和约束机制,使得获取的翻译大都是有风险的最优化搜索的结果。因此,很难说哪一种翻译模型在翻译效果上具有绝对的优势。在这样的背景下,对不同机器翻译系统的输出结果有效地进行融合以产生更好的翻译性能,即多机器翻译系统融合技术就成为机器翻译研究领域的一个不可避免的热点问题。

常用的多机器翻译系统融合可以从句子、短语和词三个级别上独立进行,同时也可以将三种方式结合起来进行融合。句子级的融合策略^[1]实际上是一种句子重排序,使用单机器翻译引擎所使用的特征以外的新特

征信息从合并的N-best列表中选择得分最高的翻译作为融合结果。句子级融合的优势在于利用了这些新的全局特征,但是这种融合不生成新的假设,而且需要估计的参数非常多,过程相对比较复杂。短语级的多系统融合方法^[1]是利用翻译候选和源语言句子的对齐重新抽取短语组成短语表,然后再进行重新解码。解码模型同标准的基于短语的机器翻译系统一样。短语级融合策略会生成新的翻译假设,扩充了翻译候选,但短语表来源于翻译候选,短语表的质量难以保证。目前国内外主要的多系统融合策略大都是词语级系统融合^[2-3],其核心思想是使用最小贝叶斯风险^[4](Minimum Bayes Risk, MBR)解码器从所有系统的N-best结果中选择一个最优结果作为骨架翻译,然后将其余翻译结果与该骨架翻译进行词对齐构建混淆网络,然后再利用投票策略重组翻译结果,产生一个得分最高的路径。该方法利用词作为单位,保证了混淆网络最大的候选空间,使得新生成的翻译假设可能性增大,但是也增加了翻译结果变差的可能性。

机器翻译模型层出不穷地出现,大大地推动了机

* 实验结果在CWMIT'2011评测会议上公开,本文在评测报告的基础上补充修改而成。本文受中国科学技术信息研究所学科建设“自然语言处理”课题(XK2011-6)、中国科学技术信息研究所重点工作“多语言信息获取关键技术研究与应用示范”课题(ZD2011-3-3)、中国科学技术信息研究所科研项目预研资金(YY-201122)和国家社科基金项目“基于本体的跨语言信息检索理论与实验研究”(06BTQ030)支持。

器翻译系统融合的发展。人们也通过评测来进一步加强系统融合的研究。我国的第四届全国机器翻译研讨会(CWMT'08)是最早开展系统融合评测项目的会议,CWMT'09继续了该项目的评测,这两次的系统融合都是针对汉英新闻领域的机器翻译任务实施融合的。遗憾的是,在之后的第七届全国机器翻译研讨会(CWMT'11)上并没有坚持系统融合的评测。知名的国际机器翻译评测NIST'09也开展了系统融合项目的评测,这也是NIST首次将系统融合作为单独的项目进行评测,但是该评测针对的任务是阿拉伯语-英语和乌尔都语-英语。

中国科学技术信息研究所(Institute of Scientific and Technical Information of China, ISTIC)作为国家级公益类科技信息研究机构,作为主要成员单位之一承担国家科技文献中心(NSTL)的建设和维护工作。掌握高端的跨语言信息技术是NSTL提供高质量信息服务的技术保障。为了更好地扩大NSTL信息检索服务的功能,ISTIC多年来一直致力于机器翻译的研究。尤其在十一五国家支撑计划项目中,针对于机器翻译模块的集成与开发这一研究目标,我们将多个机器翻译引擎集成为统一的翻译平台,作为多机器翻译系统融合的基础,进而建立了基于词和短语的多机器翻译融合系统。在2011年,ISTIC参加了第七届全国机器翻译研讨会组织机器翻译评测活动,在英汉科技领域的项目中,ISTIC采用了规则和统计两类机器翻译多引擎相结合进行系统融合的策略,提交了多机器翻译融合的结果,取得了较好的名次。本文将就ISTIC此次的参评系统作详细介绍。

本文的结构安排如下:第二节给出ISTIC机器翻译系统的总体框架和系统融合策略;第三节介绍数据的使用和处理,实验结果及相关分析在第四节;最后在第五节给出结论和展望。

2 系统描述

ISTIC提交的英汉科技领域的翻译结果为两类机器翻译多引擎的翻译输出基础上的系统融合。两类机器翻译多引擎包括基于统计的机器翻译多引擎(SMT-ME)和基于规则的机器翻译多引擎(RBMT-ME)。这两类多引擎各包含2个单引擎,每一个单引擎使用不同的参数来生成1-Best组成1-Best List,进而采用了基于词和短语的系统融合方法^[5]进行系统融合。系统的总体

框架见图1。

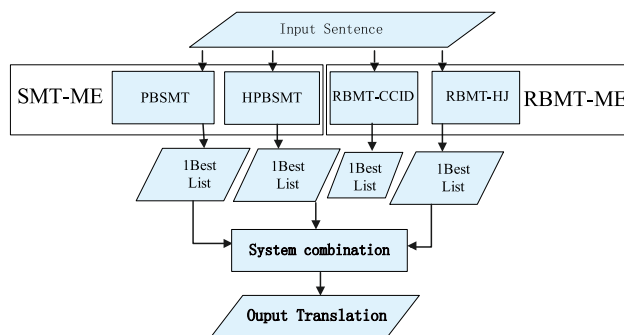


图1 ISTIC机器翻译系统框架

2.1 两类机器翻译多引擎

基于统计的机器翻译多引擎包括基于短语的统计机器翻译单引擎(Phrase Based Statistical Machine Translation, PBSMT)和基于层次短语的统计机器翻译单引擎(Hierarchical Phrase Based Statistical Machine Translation, HPBSMT)。PBSMT^[6-7]使用连续的短语对为翻译单元,能够有效地捕捉训练语料中的连续的翻译对信息。HPBSMT^[8]为基于层次短语的统计机器翻译单引擎,采用非句法知识的同步上下文无关文法,自底向上地完成目标翻译的生成。与PBSMT不同的是,该模型采用了非连续的短语,能够更为有效地对目标翻译进行排序。这两个引擎都采用了对数线性模型来进行翻译结果的遴选,采用的特征同Moses。PBSMT与HPBSMT这两个统计翻译单引擎相比较,前者是个传统的鲁棒性强的统计机器翻译单引擎,翻译原理比较简单,速度快,对大规模数据的兼容性强。后者翻译原理较为复杂,翻译速度略慢,对大规模数据的兼容性要差,但翻译效果要好。总体而言,作为统计的机器翻译,这两个单引擎的翻译效果对数据的依赖性都较强,都没有使用任何句法知识,因此翻译结果的可读性要差。

基于规则的机器翻译多引擎包含两个单引擎:赛迪规则引擎和华建规则引擎(Rule Based Machine Translation, RBMT-CCID和RBMT-HJ)。这两个引擎都是采用规则和模板相结合的技术,在传统的基于规则的机器翻译中融入了模板技术、统计技术,属于基于转换的机器翻译引擎。规则引擎对于翻译通用词汇比较擅长,其适应面较宽,翻译可读性好。但是由于科技

领域语料中专业词汇较多,这两个引擎的规则覆盖专业词汇的能力有限,因此其翻译效果要差一些。

2.2 基于词和短语的系统融合

针对两类机器翻译多引擎的特点,需要将两类翻译结果有效地进行系统融合,以期在通用词汇翻译的基础上赢得专业词汇知识的补充。

ISTIC将机器翻译的多个单引擎集成为统一的翻译平台,建立了基于词和短语的机器翻译融合系统。ISTIC采用词级的系统融合技术来构建混淆网络,将该混淆网络转换为短语表,然后使用该短语表利用短语级的系统融合技术中的重解码技术来进行解码,生成最后的融合结果。这样既保证了融合系统所构建的混淆网络的最大可能性,又可以使用更多的特征进行混淆网络解码。整个系统融合框架如图2所示。

对于每一个基于统计的单引擎,采用了不同的语言模型来生成1-Best,组成1-Best List。合并每一个单引擎的1-Best List为1-Best Lists来进行系统融合。经过在开发集上对比测试,采用了翻译效果最好的单引擎的翻译结果作为骨架翻译。在将每个翻译假设与骨架翻译进行词对齐时选用了GIZA++工具^[9]生成骨架-假设和假设-骨架双向的词对齐。GIZA++词对齐需要将骨架翻译和其余的每一个翻译假设组成平行句对,需要注意的是,这里的平行句对并不是双语的,而是单语的。由于GIZA++的词对齐质量受测试集大小的限制,为了解决这个问题,我们将所有的翻译假设中的单个词和它自身也组合成平行句对,两种平行句对合并在一起后使用GIZA++工具包进行训练,这样可以保证在词对齐的时候,相同的词肯定可以对齐。同样也采用Grow-Diag-Final式启发函数进行词对齐扩展。

利用词对齐构建混淆网络时,为了使搜索路径中的节点尽量减少不可靠的词汇,没有使用空词来扩展混淆网络。直接将对齐参考的每个词作为对照点,利用其余翻译假设与它的词对齐信息,收集在每一个翻译假设中与该词对齐的词汇作为这个词的候选翻译。这样骨架翻译的每个词会有0个或者多个候选词汇,重复的词要记录重复次数,这样就形成一个词包。将每个词包放在混淆网络的每条弧上,之后通过投票来计算每个弧上词汇的后验概率,混淆网络就构建完成。在将混淆网络转化为混淆网络短语表时,参考翻译的每个词号看作是一个短语对的源短语,它的目标短语为该词号所

对应的混淆网络上每条弧上的词,短语对的概率为该目标词的后验概率。

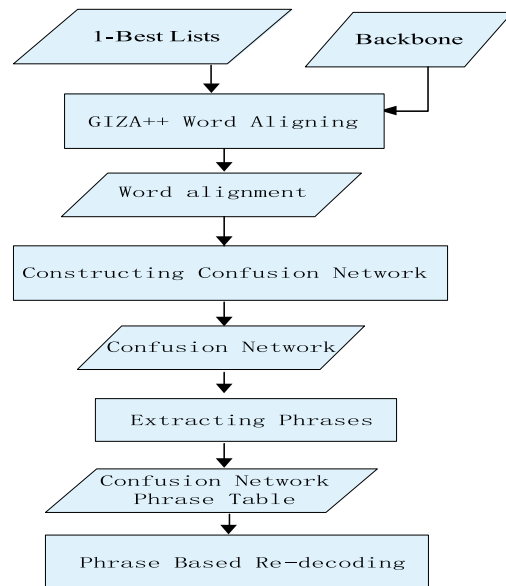


图2 ISTIC机器翻译系统融合框架

利用混淆网络短语表和一个基于短语的统计机器翻译系统来生成最后的目标翻译,这个过程类似于短语级的系统融合的再解码过程。但是这里的短语表不是源语言句子和翻译假设重新进行词对齐生成的,而是利用词级系统融合的混淆网络生成的。给定源语言句子 f ,融合的过程就是搜索具有最大概率的目标翻译 e^* :

$$e^* = \arg \max_e \sum_{m=1}^M \lambda_m h_m(e, f)$$

其中, $h_m(e, f)$ 为特征函数, λ_m 为特征权重。基于短语的重解码过程中,仍然采用对数线性模型来完成最终的翻译融合,使用的特征函数有:

- 1: 短语后验概率;
- 2: 语言模型特征;
- 3: 基于距离的重排序特征;
- 4: 词惩罚。

短语后验概率为该翻译的每个词的后验概率的对数求和,每个词的后验概率来源于生成混淆网络时该词所在词包的后验概率。语言模型特征为该翻译的每个词的语言模型概率 $p(w_n | w_{n-1} \cdots w_1)$ 的对数求和。重排序模型为基于距离的重排序模型:

$$P_{d(a_k-b_{k-1})} = |a_k - b_{k-1} - 1|$$

其中, a_k 为第 k 个目标短语的源短语的开始位置, b_{k-1} 为第 $k-1$ 个目标短语的源短语的结束位置。整句的基于距离的重排序特征为该翻译的每个短语的重排序概率对数求和。词惩罚特征为该翻译的目标词数。解码的搜索策略为柱状搜索算法, 最后生成的 **1-Best** 作为融合结果输出。每个特征函数的权重由最小错误训练算法^[10] 训练得出。

3 数据的使用和处理

训练语料采用了 CWMT' 2011 英汉科技领域项目发布的所有训练语料, 语言模型训练数据采用了训练语料的中文部分和搜狗全网新闻语料库 (SogouCA)^①, 所有的基于统计的机器翻译的单引擎的参数在该项目发布的开发集上训练。

对中文数据进行的处理有: 中文的分词和全角变半角; 对英文数据进行的处理为大小写转换和标点符号的分离处理。采用 Stanford 的中文分词^② 工具和 Moses^③ 的英文 Tokenization 工具。

词对齐工具采用了 GIZA++^④ (全部使用默认的参数) 并对该对齐结果进行扩展对齐 (Grow-Diag-Final)^[12]。

语言模型工具采用了 Srilmm^⑤ 工具包来获取 5 元文法概率信息。

4 实验

我们分别在 CWMT' 2011 英汉科技领域机器翻译项目发布的开发集和测试集上来验证翻译效果的优劣。在开发集上的打分使用了评测组织方发布的打分工具。在测试集上的打分使用了评测组织方开放的在线评测平台。

我们使用 100 个词来限制训练语料的最大长度, 获取了 896151 个平行句对作为每个统计机器翻译引擎的训练语料。表 1 列出了使用的所有语料的详细统计量。

表 1 实验语料的统计量

| 数据集 | 语言 | 句子个数 | 词汇表 | 平均句长 |
|-----|----|--------|--------|------|
| 训练集 | 中文 | 896151 | 189110 | 26 |
| | 英文 | 896154 | 180756 | 26 |
| 开发集 | 中文 | 4464 | 4658 | 21 |
| | 英文 | 1116 | 3428 | 22 |
| 测试集 | 英文 | 2497 | 6591 | 40 |

4.1 单引擎的翻译结果

表 2 列出了参与系统融合的单引擎翻译结果在开发集上的打分。其中基于短语的统计机器翻译单引擎使用了两个语言模型来进行翻译, PBSMT-1 只使用了训练语料的中文部分来训练 5 元的语言模型, PBSMT-2 在此基础上又增加了全部的 SogouCA 数据来训练 5 元的语言模型。基于层次短语的统计机器翻译单引擎使用了三个语言模型: HPBSMT-1 只使用了训练语料的中文部

表 2 单个翻译结果在开发集上的比较

| 引擎 | BLEU-SBP | BLEU | NIST |
|-----------|----------|--------|---------|
| RBMT-CCID | 0.4298 | 0.4329 | 9.6936 |
| RBMT-HJ | 0.3554 | 0.3590 | 8.7499 |
| PBSMT-1 | 0.5000 | 0.5104 | 10.5639 |
| PBSMT-2 | 0.5013 | 0.5104 | 10.6179 |
| HPBSMT-1 | 0.5061 | 0.5152 | 10.6626 |
| HPBSMT-2 | 0.5074 | 0.5177 | 10.6547 |
| HPBSMT-3 | 0.5122 | 0.5216 | 10.7158 |

分来训练 5 元的语言模型; HPBSMT-2 在此基础上增加了部分的 SogouCA 数据 (大约 200 万句); HPBSMT-3 在此基础上增加了全部的 SogouCA 数据。从表 2 可以看出, 在英汉科技领域, 基于统计的翻译多引擎的翻译表

^① <http://www.sogou.com/labs/dl/ca.html>.

^② <http://www-nlp.stanford.edu/downloads/segmenter.shtml>.

^③ <http://www.statmt.org/ Moses/>.

^④ <http://giza-pp.googlecode.com/>.

^⑤ <http://www.speech.sri.com/projects/srilmm/download.html>.

现要优于基于规则的翻译多引擎,基于层次短语的翻译单引擎的表现要优于基于短语的翻译单引擎。不同的语言模型也给出了不同的翻译结果,语言模型数据用得越多,翻译结果越好,但是翻译效果增长的幅度并不大。

4.2 系统融合结果

除了基于词和短语的系统融合(WPSC)方法,我们也尝试了其他系统融合方法,包括TER^[1,4]、WER^[13]、INHMM^[14]。WPSC选用了HPBSMT-3的翻译结果作为骨架翻译。表3列出了系统融合的打分。其中RULE表示有规则结果参与系统融合,即使用了表2中的7个翻译结果。反之,NO-RULE表示没有规则结果参与系统融合,即只使用了表2中的从第4行到第8行的5个翻译结果。从表2的结果来看,在英汉科技文献翻译中,基于规则的机器翻译多引擎效果不如基于统计的机器翻译多引擎,但从表3的结果来看,使用规则引擎的翻译结果和不使用规则引擎的翻译结果来比较,前者融合效果普遍比后者稍好,可见基于规则的机器翻译多引擎能在一定程度上弥补基于统计机器翻译多引擎的一些不足,提高系统融合的翻译结果的质量。在表3中,我们使用的WPSC与其他融合方法的对比来看,WPSC的效果比TER、WER融合方法略好;与INHMM方法比较,使用规则系统融合的结果稍差;但不使用规则系统融合的结果则略好。由此可见,WPSC融合方法是一种既直观又有效的融合方法。

我们也在此次机器翻译评测的英汉科技领域的测试集上进行了实验。表4列出了单引擎和系统融合在测试集上的结果。分析实验结果说明,基于词和短语的系统融合方法(WPSC)能够达到基本的融合要求,能保证融合结果略优于最好的翻译引擎的翻译结果,但是融合的效果没有超过INHMM,这个结果与在开发集上的实验是一致的。因此,我们在评测中提交了INHMM的融合作为primary结果,WPSC的结果作为contrast结果。在最后的公布的评测结果中,这两个融合结果都取得了不错的成绩。

WPSC是个比较保守的系统融合方法。其基本思想在于试图在最好的翻译结果的基础上使用其他翻译假设的词汇来进行补充,以达到比参与融合的最好单个翻译结果更好。

表3 系统融合在开发集上的翻译结果比较

| 数据集 | BLEU-SBP | BLEU | NIST | |
|-------|----------|--------|--------|---------|
| TER | RULE | 0.5127 | 0.5240 | 10.6997 |
| | NO-RULE | 0.5070 | 0.5181 | 10.6403 |
| WER | RULE | 0.5113 | 0.5205 | 10.6412 |
| | NO-RULE | 0.5047 | 0.5183 | 10.5810 |
| INHMM | RULE | 0.5268 | 0.5438 | 10.7900 |
| | NO-RULE | 0.5084 | 0.5250 | 10.6018 |
| WPSC | RULE | 0.5156 | 0.5248 | 10.7829 |
| | NO-RULE | 0.5131 | 0.5221 | 10.7436 |

5 讨论

目前的机器翻译系统模型众多,但是真正一枝独秀的模型很少。机器翻译研究如果想获得大规模应用性发展,系统融合研究是必不可少的。系统融合不但可以博采众家机器翻译模型在理论研究的长处,即使对于同一翻译模型内部,不同翻译参数输出的翻译结果也具备一定的融合空间,更甚之,对于不同的数据集生成的翻译结果也可以一并融合之。

目前系统融合可以在句子、短语和词三个级别上进行,关键技术主要集中在词级系统融合的词对齐研究上。系统融合效果的好坏可以取决于下列因素:

(1) 参与融合的单系统的翻译结果的选择:系统融合效果的好坏很大程度上取决于参与融合的单系统的翻译结果。无论是翻译模型的差异性,还是同一个翻译模型采取不同的参数,亦或是N-Best数量的选择,都会影响系统融合的质量。到底什么样的融合策略具备足够鲁棒的翻译性能,目前还没有定论。评测中常常是在开发集上尝试过不同的组合策略之后再选取最有效的方式。我们的经验是:单系统的翻译结果在翻译模型上最好既有差异性又能够互补,翻译质量不能相差太大,否则差的单个翻译结果会导致系统融合的翻译结果比不上最优的单个系统的翻译结果。

(2) 系统融合策略的制定:系统融合的目的是获取在现有单个系统的翻译结果基础上更为优质的翻译,但是目前系统融合的方法都没有办法做到足够的鲁棒,常常是在不同的数据集上有不同的表现,甚至不能超过最优的单个系统的翻译结果。因此,如何制

表4 系统融合在测试集上的翻译结果比较

| Results | BLEU5-SBP | BLEU5 | BLEU6 | NIST6 | NIST7 | GTM | mWER | mPER | ICT |
|------------|-----------|--------|--------|---------|---------|--------|--------|--------|--------|
| RBMT-CCID | 0.3079 | 0.3150 | 0.2551 | 9.0899 | 9.1040 | 0.8099 | 0.6151 | 0.3651 | 0.2819 |
| RBMT-HJ | 0.2408 | 0.2476 | 0.1930 | 8.1135 | 8.1223 | 0.7588 | 0.6739 | 0.4144 | 0.2364 |
| PBSMT-1 | 0.3860 | 0.3989 | 0.3321 | 10.4192 | 10.4415 | 0.8388 | 0.6147 | 0.3085 | 0.3852 |
| PBSMT-2 | 0.3881 | 0.4020 | 0.3353 | 10.4731 | 10.4947 | 0.8342 | 0.6140 | 0.3038 | 0.3927 |
| HPBSMT-1 | 0.3853 | 0.3970 | 0.3308 | 10.4071 | 10.4307 | 0.8380 | 0.6064 | 0.3080 | 0.3808 |
| HPBSMT-2 | 0.3986 | 0.4120 | 0.3463 | 10.5815 | 10.6075 | 0.8395 | 0.6060 | 0.3036 | 0.3936 |
| HPBSMT-3 | 0.3965 | 0.4102 | 0.3438 | 10.5567 | 10.5823 | 0.8395 | 0.6132 | 0.3047 | 0.3917 |
| INHMM-RULE | 0.4083 | 0.4302 | 0.3656 | 10.3527 | 10.3794 | 0.8219 | 0.5687 | 0.3206 | 0.4264 |
| WPSC-RULE | 0.4009 | 0.4142 | 0.348 | 10.6314 | 10.6573 | 0.8435 | 0.6041 | 0.3004 | 0.3979 |

定系统融合策略,使得即使得不到更好的翻译候选的组合,也能选出最优的单个系统的翻译结果,这是个值得努力的方向。我们的系统融合方法在这方面进行了尝试,并选择了一个保守的融合策略,通过固定最优的单一系统的翻译结果作为骨架翻译,使用其他翻译结果中的词汇来替换骨架翻译中的词汇来生成目标翻译,这个方法目前只使用了简单的重排序模型和语言模型来保证目标翻译的质量,翻译表现还有待提高。

(3) 翻译效果的评价:系统融合策略的制定常常是根据翻译效果的评价来制定,因此翻译效果的评价起到了至为关键的作用。机器翻译评价包括人工评价和自动评价两种,前者代价昂贵,因此主要使用了BLEU打分等自动评价。自动评价不仅体现在最终系统融合效果的评价,还表现在系统融合过程中骨架翻译的选择。而BLEU值还存在很大的局限性,并不能完全有效地去指导机器完成这些任务。因此,如何制定算法来捕捉参与融合的单个系统的翻译结果中对目标翻译有贡献的片段是个值得研究的方向。系统融合的发展目前还处于词和短语阶段。通过源语言或目标语言的句法和语义知识来深层次地指导融合,将能较好地克服系统融合中目前所困扰的译文不连续或译文不符合

语法结构、融合性能不稳定等难题,最终得到多种翻译方法的水乳交融^[15]。

另外,选取不同的机器翻译引擎来获取更好的融合结果的时候,也要考虑到在具体应用中,由于多个引擎同时翻译后才能融合,因此需要更长的时间来等待。如何并发翻译、快速融合是多机器翻译系统融合发展的一个趋势。

6 结论与展望

ISTIC参加了CWMT2011英汉科技领域评测任务的机器翻译项目,并取得了较好的成绩。

囿于翻译的效率问题,目前,我们所采用的融合系统的技术还比较简单,仍有大量的空间可以进一步加强研究。我们的参考翻译策略比较简单,对于所有的句子都采用了单一系统的结果做参考,可以采用最小贝叶斯风险解码技术或者其他评价机制来动态地选择参考翻译。我们的词对齐策略也比较简单,只使用了GIZA++工具来生成,可以采用更多的启发式措施来克服未登录词的问题。

参考文献

- [1] ROSTIA, AYAN N, XIANG B, et al. Combining outputs from multiple machine translation systems [C]// Proceedings of NAACL-HLT, 2007: 228-235.
- [2] FISCUS J. A post-processing system to yield reduced word error rates: recognizer output voting error reduction (ROVER) [C]// IEEE Workshop on Automatic Speech Recognition and Understanding, 1997: 347-354.
- [3] SCHWENK H, GAUVAIN J. Improved ROVER using language model information [C]// ISCA ITRW Workshop on Automatic Speech Recognition: Challenges for the new Millennium, 2000: 47-52.

- [4] SIM K, BYRNE W, GALES M, et al. Consensus network decoding for statistical machine translation system [C]// IEEE International Conference on Acoustics, Speech, and Signal Processing, 2007: 105-108.
- [5] 何彦青, 张均胜, 王惠临. 基于词和短语的多机器翻译系统融合方法研究[J]. 情报学报, 2011(12).
- [6] KOEHN P, OCH F, MARCU D. Statistical phrase-based translation [C]// Proceedings of the Human Language Technology conference / North American Chapter of the Association for Computational Linguistics (HLT/NAACL-2003), 2003: 127-133.
- [7] OCH F, NEY H. The Alignment template approach to statistical machine translation [J]. Computational Linguistics, 2004(30): 417-449.
- [8] CHIANG D. A hierarchical phrase-based model for statistical machine translation [C]// Proceedings of ACL 2005, 2005: 263-270.
- [9] OCH F, NEY H. A systematic comparison of various statistical alignment models [J]. Computational Linguistics, 2003(29): 19-51.
- [10] OCH F. Minimum error rate training in statistical machine translation [C]// Proceedings of the 41st Annual Meeting on Association for Computational Linguistics, Sapporo, Japan, 2003.
- [11] KOEHN P. Pharaoh: A beam search decoder for phrase-based statistical machine translation models [C]// Proceedings of the 6th Conference of the Association for Machine Translation in the Americas, 2003: 115-124.
- [12] STOLCKE A. SRILM-an extensible language modeling toolkit [C]// Proceedings of International Conference on spoken language processing, 2002: 901-904.
- [13] BANGALORE S, BORDEL G, RICCARDI G. Computing consensus translation from multiple machine translation systems [C]// Proceedings of ASRU, 2001: 351-354.
- [14] HE X, YANG M, GAO J, et al. Indirect-HMM-based hypothesis alignment for combining outputs from machine translation systems [C]// Proceedings of EMNLP, 2008.
- [15] 李茂西, 宗成庆. 机器翻译系统融合技术综述[J]. 中文信息学报, 2010(4): 74-84.

作者简介

何彦青, 博士, 模式识别与智能系统专业, 主要研究方向: 机器翻译, 自然语言处理。E-mail: heyq@istic.ac.cn
石崇德, 博士, 情报学专业, 主要研究方向: 机器翻译, 中文分词。E-mail: shicd@istic.ac.cn
于薇, 硕士, 情报学专业, 主要研究方向: 自然语言处理, 电子文件管理。E-mail: yuwei@istic.ac.cn
张均胜, 博士, 计算机软件与理论专业, 主要研究方向: 多语言信息服务, 语义计算。E-mail: zhangjs@istic.ac.cn
王惠临, 研究员, 博士生导师, 主要研究方向: 多语言信息服务, 机器翻译, 自然语言处理。E-mail: wanghl@istic.ac.cn

Machine Translation System Combination and Its Application

He Yanqing, Shi Chongde, Yu Wei, Zhang Junsheng, Wang Huilin / Information Technology Support Center, Institute of Scientific and Technical Information of China, Beijing, 100038

Abstract: Machine translation system combination has become a hot topic in machine translation research in recent years, which combines the outputs of different machine translation systems to get a better translation result. This paper is an overview of ISTIC evaluation technical report for the 7th China workshop on machine translation. ISTIC participated in the English-to-Chinese machine translation task in scientific and technical domain. This paper describes the implement framework of our machine translation system. We also discuss the key techniques and analyze the experimental results on the evaluation data. Finally, we discuss some influence factors on system combination. The future development prospects of our system combination are also discussed.

Keywords: Machine translation, Natural language processing, System combination

(收稿日期: 2011-11-03)