

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5705472号
(P5705472)

(45) 発行日 平成27年4月22日 (2015. 4. 22)

(24) 登録日 平成27年3月6日 (2015. 3. 6)

(51) Int. Cl.

F I

G 0 6 F 17/28 (2006.01)

G 0 6 F 17/28 6 1 8

請求項の数 3 (全 21 頁)

(21) 出願番号	特願2010-166695 (P2010-166695)	(73) 特許権者	596170170
(22) 出願日	平成22年7月26日 (2010. 7. 26)		ゼロックス コーポレイション
(65) 公開番号	特開2011-28754 (P2011-28754A)		XEROX CORPORATION
(43) 公開日	平成23年2月10日 (2011. 2. 10)		アメリカ合衆国、コネチカット州 068
審査請求日	平成25年7月23日 (2013. 7. 23)		56、ノーウォーク、ビーオーボックス
(31) 優先権主張番号	12/509, 633		4505、グローバー・アヴェニュー 4
(32) 優先日	平成21年7月27日 (2009. 7. 27)		5
(33) 優先権主張国	米国 (US)	(74) 代理人	110001210
			特許業務法人 Y K I 国際特許事務所
		(72) 発明者	ミハイル ザスラフスキ
			フランス ジフ シュ リヴェット リュ
			デュ クロス 1
		(72) 発明者	マルク ディムトマン
			フランス グルノーブル キ ド フラン
			ス 4
			最終頁に続く

(54) 【発明の名称】 一般化された巡回セールスマン問題としてのフレーズベースの統計的機械翻訳

(57) 【特許請求の範囲】

【請求項 1】

統計的機械翻訳 (SMT) および一般化された非対称巡回セールスマン問題 (GTSP)
) グラフを使用してソース言語を目標言語に翻訳する方法であって、
 コンピュータが、

ソース言語及び目標言語で表現されたフレーズのペアであるバイ - フレーズを各ノード
 とし、前記各ノード間のエッジにコストが設定された GTSP を形成するステップと、
 ソース言語で表現された入力文に含まれる各フレーズを、前記 GTSP を表す GTSP
 グラフ内のノードに対応するバイ - フレーズに基づいて、目標言語に置換するステップと、

前記入力文のブロックを目標言語での表現に置換するために用いた前記ノードを巡回する
 巡回経路であって、コストの合計が最小となる巡回経路を決定するステップと、

前記巡回経路によって定義される順序で目標言語に置換された各フレーズを出力するス
 テップと、

を実行し、

前記巡回経路を決定するステップは、

GTSP グラフの最適巡回を生成するステップと、

前記 GTSP グラフの 3 つ以上のノードの接続関係に基づいて決定されたコストに基づ
 いて、前記最適巡回の真のコスト C_{+} を計算するステップと、

前記 GTSP グラフのエッジに設定されたコストに基づいて、前記最適巡回の見かけの

10

20

コスト C_a を計算するステップと、

前記真のコスト C_t と前記見かけのコスト C_a の間の差 D を判定するステップと、

前記差 D があらかじめ設定された閾値 より小さいか否かを判定するステップと、

D が前記あらかじめ設定された閾値 より小さい場合には、前記巡回経路として前記最適巡回を出力するステップと、

を包含することを特徴とする方法。

【請求項 2】

前記巡回経路を決定するステップは、

前記 $G T S P$ グラフを非対称巡回セールスマン問題 ($A T S P$) グラフに変換するステップと、

前記 $A T S P$ グラフを標準巡回セールスマン問題 ($T S P$) グラフに変換するステップと、

前記 $T S P$ グラフに対してコンコード ($Concorde$) ソルバおよびリン・カーニハンのヒューリスティックのうちの少なくとも 1 つを使用してコストの合計が最小となる巡回経路を決定するステップと、

を包含する請求項 1 に記載の方法。

【請求項 3】

さらに、

D が前記あらかじめ設定された閾値 に等しいか、またはそれを超える場合に、前記グラフ内の少なくとも 1 つのノードを絞り込み、前記絞り込んだノードを包含する絞り込み済みグラフを生成するステップと、

D が より小さくなるまで反復的に、1 つまたは複数の絞り込み済みグラフについて前記真のコスト C_t および見かけのコスト C_a を計算し、それらの間の前記差 D を判定し、かつ前記差 D と前記あらかじめ設定された閾値 を比較するステップと、を包含し、

少なくとも、

前記最適巡回が各ノードを正確に一度だけ訪問し、前記最適巡回内のノードの間の各エッジがそれぞれのバイグラム重みと関連付けられ、かつ前記最適巡回の前記見かけのコスト C_a が前記巡回内のすべてのエッジの前記バイグラム重みの合計によって計算されることと、

前記最適巡回の前記真のコスト C_t がトリグラム・コストを使用して計算されることと、のうちの一方を含む、

請求項 1 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

この出願は、コンピューティング・システム内における統計的機械翻訳 ($statistical\ machine\ translation, SMT$) に関する。ここで述べる手法は、このほかの翻訳システム、このほかの統計的マッピング応用、および/またはこのほかの翻訳方法の中にも応用を見つけることができることを理解されたい。

【背景技術】

【0002】

統計的機械翻訳 (SMT) に対する古典的アプローチは、「バイ・フレーズ ($bi-phrase$)」を伴う。「バイ・フレーズ」とは、ソース言語および目標言語の表現またはフレーズのペアであり、この表現またはフレーズのペアは、ソース文から目標 (すなわち、翻訳された) 文を構成するためのビルディング・ブロックを形成する。

【0003】

N -グラム ($N-gram$) 言語モデルは、シーケンス内の次の項目を予測するための確率論的モデルの一種である。 N -グラムは、統計的自然言語処理および遺伝子配列の分析の多様な分野において使用されている。 N -グラムは、与えられたシーケンスからの n 個の項目のサブシーケンスである。懸案の項目は、音素、音節、文字、単語、塩基対等々

10

20

30

40

50

とすることができる。

【 0 0 0 4 】

与えられたソース文 S を翻訳するために、古典的なフレーズ・ベースの SMT システムは、次の形式の対数 - 線形モデル (\log - $l i n e a r$ $m o d e l$) を用いる。

【 数 1 】

$$p(t, a | s) = 1/Z_s \exp \sum_k \lambda_k h_k(s, a, t)$$

ここで、 h_k は、「特徴」であり、ソース文字列 s 、目標文字列 t 、および配列 a の関数である。配列 a は、ソース文字列 s から目標文字列 t を組み立てるのに用いられる、パイ - フレーズのシーケンスの表現である。 λ_k は、重みであり、 Z_s は、 $p(t, a | s)$ がペア (t, a) について適正な条件付き確率分布となることを保証する正規化因子である。

10

【 0 0 0 5 】

一旦、対数 - 線形モデルが定義されてしまえば (トレーニング段階を伴う ; たとえば、参照によりこれに援用される非特許文献 1 を参照されたい)、デコードの役割は、条件付き確率 $p(t, a | s)$ を最大化するペア (t, a) を見つけること、および対応する目標文字列 t を出力することになる。

【 0 0 0 6 】

古典的なシステムは、ヒューリスティックな左から右へのサーチの何らかの変形に基づいており、これは、各ステップにおいて、新しいパイ - フレーズを用いて現在の部分翻訳を拡張しつつ、かつ 2 つのスコア、すなわちこれまでの部分翻訳の既知の要素についてのスコア、および翻訳を完了するための残りのコストのヒューリスティックな評価を計算しつつ、左から右へと漸増的に候補翻訳 (t, a) の組み立てを試みる。もっとも頻繁に使用される変形は、いくつかの部分的な候補が並列に維持され、現在の評価が低すぎる候補が取り除かれて、より有望な候補を選ぶビーム - サーチの形式である。

20

【 先行技術文献 】

【 非特許文献 】

【 0 0 0 7 】

【 非特許文献 1 】 ロペズ・A (Lopez, A.) 著、2008. Statistical Machine Translation. ACM Comput. Surv. 40, 3 (2008 年 8 月)、p. 1 - 49

30

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 8 】

左から右へのヒューリスティックなサーチによる翻訳を行う従来技術では、サーチの早い段階で生じる誤った選択から過剰な影響を受けることがあり、早い段階での誤りから回復することは困難であり得る。

【 課題を解決するための手段 】

【 0 0 0 9 】

グラフ内のノードとしてパイ - フレーズをモデリングすることによってバイグラム (または、より高い N - グラム) 言語モデルを組み込むフレーズ・ベースのモデルのための SMT を容易にするシステムおよび方法を述べる。たとえば、統計的機械翻訳 (SMT) および一般非対称巡回セールスマン問題 (generalized asymmetric traveling salesman problem, GTSP) のグラフを使用して 2 つの言語を翻訳する方法は、GTSP として SMT を定義するステップと、GTSP を表す GTSP グラフ内のノードに対応するパイ - フレーズを使用して入力文のブロックを翻訳するステップと、GTSP を解決するステップと、GTSP 解決によって定義される順序で翻訳されたブロックを出力するステップと、を包含する。

40

【 図面の簡単な説明 】

50

【 0 0 1 0 】

【図 1】 G T S P アプローチを使用して言語間の翻訳において使用するためのグラフを生成する、フレーズ - ベースの S M T を実行するためのシステムを図解したブロック図である。

【図 2】それぞれの「コスト」に従って 0 ~ 6 がラベル付けされた複数のエッジを伴う非対称 T S P (a s y m m e t r i c T S P , A T S P) を標準 T S P に変換するための第 1 の変換を図解した説明図である。

【図 3】 A T S P を標準 T S P に変換するための第 2 の変換を図解した説明図である。

【図 4】 G T S P を A T S P に変換する変換を図解した説明図である。

【図 5】既存のエッジのサブセット、すなわちノード「 t r a d u c t i o n - m t 」に入るか、または出るすべてのエッジだけが示された、ソース文「 c e t t e t r a d u c t i o n a u t o m a t i q u e e s t c u r i e u s e 」についての遷移グラフを図解した説明図である。

【図 6 A】 1 つの出力に対応する G T S P 巡回を図解した説明図である。

【図 6 B】他の 1 つの出力に対応する G T S P 巡回を図解した説明図である。

【図 7】パイ - フレーズ「 i 」だけが取り除かれたグラフであって、現在は「 i 」をカプセル化している拡張パイ - フレーズのいくつかだけがグラフを通る 1 つの有効な巡回を定義するエッジを伴って示されているグラフを図解した説明図である。

【図 8】選択的なオンデマンド絞り込みを図解した説明図である。

【図 9 A】 T S P としてフレーズ - ベースの統計的機械翻訳を実行するための方法を図解したフローチャートである。

【図 9 B】 T S P としてフレーズ - ベースの統計的機械翻訳を実行するための代替方法を図解したフローチャートである。

【図 1 0】トリグラム言語モデルを使用してフレーズ - ベースの翻訳を実行するためのグラフを図解した説明図である。

【図 1 1】差が より大きい、または等しい場合に得られるグラフを図解した説明図である。

【発明を実施するための形態】

【 0 0 1 1 】

グラフ内のノードとしてパイ - フレーズをモデリングすることによってフレーズ - ベースの S M T を容易にするシステムおよび方法を述べる。ここで使用するときの「グラム」は、単語を意味するものであり、バイグラム言語モデルは 2 語のグループを採用し、トリグラム言語モデルは 3 語のグループを採用し、以下同様とする。

【 0 0 1 2 】

このシステムおよび方法においては、パイ - フレーズが、グラフ内のノードとしてモデル化される。それに加えて翻訳の構成が、グラフ内のノードの間の「巡回」、すなわち各ノードを正確に一度だけ訪問する経路としてモデル化される。巡回の全体的なコストが、その巡回の間に通り抜けたエッジに関連付けされたコストを加算することによって計算される。

【 0 0 1 3 】

したがって、ここで述べるシステムおよび方法は、 S M T 問題を G T S P 問題に直接マップし、 G T S P としてフレーズ - ベースの翻訳を表現する。

【 0 0 1 4 】

図 1 を参照すると、 G T S P アプローチを使用してフレーズ - ベースの S M T の実行を容易にするシステム 1 0 が図解されている。このシステムは、ここで述べる多様な手法、方法、応用、アルゴリズム等を実行するためのコンピュータ実行可能命令を実行するプロセッサ 1 2 およびそれらの命令を記憶するメモリ 1 3 を含む。

【 0 0 1 5 】

メモリ 1 3 は、コンピュータ上で実行できるコンピュータ・プログラム製品を包含し得る。コンピュータ・プログラム製品は、コントロール・プログラムが記録されたコンピュ

10

20

30

40

50

ータ可読記録媒体（たとえばメモリ 13）であってよく、例えば、ディスク、ハード・ドライブ、またはこれらと同種のものであってよい。一例によれば、ここで述べる手法は、1つまたは複数の汎用コンピュータ、専用コンピュータ（1つまたは複数）、プログラムされたマイクロプロセッサまたはマイクロコントローラおよび周辺集積回路素子、ASIC またはそのほかの集積回路、デジタル信号プロセッサ、ハードワイヤードの電子または論理回路、たとえばディスクリート素子回路、プログラマブル論理デバイス、たとえば PLD、PLA、FPGA、グラフィック・カード CPU（GPU）、または PAL、またはこれらの類の上で実装することができる。

【0016】

システム 10 は、さらに、ユーザがシステムによる翻訳のための入力文 15 を入力することができるユーザ・インターフェース 14 を包含している。入力文 15 は、コンピュータ実行可能アルゴリズム（1つまたは複数）を使用してプロセッサ 12 によって処理され、オプションとして、翻訳済みデータ 17（たとえば、翻訳された文）がユーザ・インターフェース 14 に出力されるまで、1つまたは複数の中間データ段階を通過する。

10

【0017】

メモリ 13 は、さらに、入力文 15 を翻訳するための多様な構成要素（たとえば、コンピュータ実行可能インストラクションまたはその類）を記憶する。それに加えて、あらかじめ生成済みのバイ・フレーズがバイ・フレーズ・ライブラリ 20 内に記憶される。

【0018】

翻訳のための入力文を受け取ると、プロセッサ 12 は、入力文と両立するバイ・フレーズをバイ・フレーズ・ライブラリ 20 から検索するとともに、言語モデル 19 にもアクセスし、検索したバイ・フレーズおよびその言語モデルを利用して GTS P グラフ 22 を組み立てる。

20

【0019】

GTS P グラフ 22 のノードを通る最適巡回を生成するために、プロセッサは、厳密ソルバ・アルゴリズム 24、近似ソルバ・アルゴリズム 25 等々のうちの 1つまたは複数であり得る TSP ソルバ 23 を実行する。

【0020】

システム 10 は、GTS P グラフ 22 を入力文 15 に適用することによって入力文のフレーズ・ベースの統計的機械翻訳（phrase based statistical machine translation, PB SMT）を容易にする。プロセッサ 12 は、GTS P として PB SMT タスクを定義する。プロセッサは、入力文と矛盾のない 1つまたは複数のバイ・フレーズをバイ・フレーズ・ライブラリ 20 から検索し、それぞれがバイ・フレーズに対応する複数のノードを包含する GTS P グラフ 22 を生成する。TSP ソルバ 23 が実行され、GTS P グラフ 22 の最適巡回が生成される。

30

【0021】

1つの実施態様においては、バイグラム言語モデルに代えて N - グラム言語モデル（2より大きい N を用いる）を使用するとき、プロセッサは、最適巡回の真のコスト C_t 、最適巡回の見かけのコスト C_a を計算し、真のコスト C_t と見かけのコスト C_a の間の差 D を決定する。プロセッサは、D があらかじめ決定済みの閾値 より小さいとき、GTS P に対する解として最適巡回を出力し、出力される GTS P 解を使用して入力文を第 1 の言語から第 2 の言語へ翻訳する。

40

【0022】

プロセッサ 12 は、D があらかじめ決定済みの閾値 に等しいか、またはそれを超える場合に、GTS P グラフ 22 内の少なくとも 1つのノードを絞り込み、絞り込んだノードを包含する絞り込み後のグラフを生成する。プロセッサは、グラフの各絞り込みについて真のコスト C_t および見かけのコスト C_a の計算、それらの間の差 D の決定、および差 D とあらかじめ決定済みの閾値 の比較を、D が より小さくなるまで反復的に継続する。

【0023】

TSP モデルは、次に述べるとおり、4つの主要な変形を含む。対称 TSP (symm

50

etric TSP, STSP) は、 N 個のノードについての無向グラフ G を伴い、このグラフは、エッジ (ライン) が実数値コストを持ち、 $+$ (正の無限大) のコストが許容される。STSP 問題は、合計コストが最小となる「巡回」を見つけ出すことにあり、それにおいて巡回 (ハミルトン閉路 (Hamiltonian Circuit) と呼ばれる) は、グラフの各ノードを正確に一度だけ訪問するノード $X_1, X_2, \dots, X_N, X_1$ の「循環」シーケンスであり、巡回の合計コストは、対応するエッジの寄与を加算することによって計算される。

【0024】

ATSP は STSP の変形であり、基礎をなすグラフ G が有向であり、グラフの 2 つのノード a および b について、エッジ (a, b) とエッジ (b, a) とが異なるコストを持

10

【0025】

一般対称 TSP (Generalized Symmetric TSP) または SGTSP は、エッジが実数値コストを持つ N 個のノードの無向グラフ G を伴う。 N 個のノードを、 M 個の空でない共通の要素を持たないサブセット (クラスタと呼ばれる) へ分割することを考えると、目的は、各クラスタが正確に一度だけ訪問される最小合計コストを伴う M 個のノード $X_1, X_2, \dots, X_M, X_1$ の循環シーケンスを見つけることとなる。

【0026】

一般非対称 TSP (Generalized Asymmetric TSP) または GTSP は、SGTSP に類似であるが、グラフ G が有向グラフである。理解されるものとするが、ここで GTSP を使用して説明する場合には、特に示さない限りは非対称 GTSP を意味する。

20

【0027】

STSP は、しばしば単に TSP と示され、NP 困難であることが知られているが、そのための効率的な厳密および近似ソルバの開発には多大な関心が存在し、それにおいては「効率」が、いわゆる TSP LIB ライブラリに提供されるような大規模ベンチマーク例を解決するために要する時間によって測定される。

【0028】

ATSP、SGTSP、および GTSP は、すべて、STSP への単純な (たとえば、問題のインスタンスのサイズにおける多項式または線形増加) 変換によってマップが可能である。たとえば以下に述べるとおり、2 つの「変換器」(たとえば、メモリ内に記憶され、プロセッサによって実行されるプログラム) が採用される。このうち 1 つは GTSP から ATSP への変換のための変換器であり、もう 1 つは ATSP から TSP への変換のための変換器である。

30

【0029】

引き続き図 1 を参照するが、図 2 に、それぞれの「コスト」に従って 0 ~ 6 がラベル付けされた複数のエッジを伴う第 1 の変換 30 を図解する。たとえばアプリゲート (Applegate) ほかの「The Traveling Salesman Problem: A Computational Study」Princeton UP、2006 年、p. 126 を参照されたい。元の有向グラフ 32 の各ノード A は、変換された無向グラフ 34 の 3 つのノード A, A', A'' によって置換され、2 つの「0 コスト」のエッジが追加される。無向グラフの任意の巡回中に中間ノード A' が存在しなければならないことから、0 コストのエッジも存在しなければならない。またこのことが、たとえば、最初にコスト 2 のエッジを通り、続いてコスト 5 のエッジ (および、同様に「逆」方向に対応することになるあらゆるエッジのペア) を通る巡回を除外する。これは、その後その巡回がノード A に入射する 3 つのエッジを含む必要があり、それは不可能であるからである。したがって、無向グラフ 34 の任意の最適巡回は、元のグラフ 32 の最適巡回に対応する。

40

【0030】

50

図3は、いくつかのエッジへの大きな人工的な重みの導入を対価として、2つのノードを導入するだけで元のグラフの1つのノードを置換する利点を有する第2の変換50を図解している。元のグラフ52の各ノード X （たとえば、ノードA、B、およびC）が、ノード X および X' に複製され、大きな負の重み $-K$ で X と X' とを接続し、元のグラフ内の有向エッジ (X, Y) のコストが変換後のグラフ54内のエッジ (X', Y) 上に再現される。 K が十分に大きい（たとえば、元のグラフ52内のすべての有限のコストの合計より K が大きい）場合には、無向グラフ54内の任意の最適巡回が、ほかのいずれの構成より (X, X') エッジを通ることを選択する。これらの (X, X') エッジのうち1つでも通らないことがあれば、少なくとも K 単位分のコストを失うことを意味するからである。しかしながらこれは、変換後のグラフ54のノードの任意の最適巡回において、 X と X' とが常に互いに隣り合い、 X と Y との間または X' と Y' との間にリンクが存在しないことから、その種の巡回だけが $X_1, X'_1, X_2, X'_2, \dots, X_N, X'_N, X_1$ 、または $X'_1, X_2, X'_2, \dots, X_N, X'_N, X_1, X'_1$ という形式になることを意味し、これは、元のグラフ52における「方向の変更」を禁止する制約条件に対応する。

10

【0031】

図4は、GTSPをATSPに変換する変換70を図解している。たとえば、ヌーン(Noon)ほかの「An efficient transformation of the generalized traveling salesman problem」INFOR31(1993年)p.39~44を参照されたい。この変換においては、元のグラフ74内の Y_1, \dots, Y_k が与えられたクラス72のノードであり、 X および Z がほかのクラスに属する任意のノードであると仮定する。変換後のグラフ76においては、図に示されるとおりに循環を形成するために Y_i の間にエッジ78が導入され、各エッジは大きな負のコスト $-K$ を有する。 X から Y_i へ入るエッジはそのまま残され、 Y_i から Z へ出るエッジがその起点を Y_{i-1} に変更されている。すると、 X, Y_i, Z を通過する元のGTSP問題における実行可能な巡回は、変換後のグラフ76において最初に X を通り、続いて $Y_i, Y_{i+1}, \dots, Y_k, \dots, Y_{i-1}$ を通り、その後 Z を通る巡回として「エンコード」される（このエンコードは元のコストから $(k-1)K$ を減じたコストを有する）。それに加えて K が十分に大きければ、変換後のATSPグラフのためのソルバが、可能な限り多くの $-K$ エッジを通り抜ける傾向を有することになり、これは、正確に $k-1$ 個のその種のエッジ、たとえばそのクラスに関連付けされた1つのエッジを除くすべてを通り抜けることを意味する（その種のエッジすべてを通り抜けるとグラフ全体のための巡回を見つけることができなくなるため、ソルバがそれを行うことはない）。言い替えると、GTSP問題のいくつかの実行可能な巡回のエンコーディングである巡回を作り出すことになる。

20

30

【0032】

以下の例は、巡回セールスマン問題としてフレーズ-ベースのデコーディングを説明する。この例では、フランス語の文「cette traduction automatique est curieuse（この機械翻訳は奇妙である。）」が英語に翻訳される。この文を翻訳のための関連するパイ-フレーズを次の表1に示す。

40

【表 1】

バイーフレーズ 識別子	ソース	目標
h	<i>cette</i>	<i>this</i>
t	<i>traduction</i>	<i>translation</i>
ht	<i>cette traduction</i>	<i>this translation</i>
mt	<i>traduction automatique</i>	<i>machine translation</i>
a	<i>automatique</i>	<i>automatic</i>
m	<i>automatique</i>	<i>machine</i>
i	<i>est</i>	<i>is</i>
s	<i>curieuse</i>	<i>strange</i>
c	<i>curieuse</i>	<i>curious</i>

10

【0033】

このモデルにより、次の翻訳が生成される。

h . m t . i . s t h i s m a c h i n e t r a n s l a t i o n i s s t r a n g e (この機械翻訳は奇妙である)

20

h . c . t . i . a t h i s c u r i o u s t r a n s l a t i o n i s a u t o m a t i c (この好奇心旺盛な翻訳は自動である)

h t . s . i . a t h i s t r a n s l a t i o n s t r a n g e i s a u t o m a t i c (この翻訳奇妙は自動である)

...

上記では、各翻訳を導くバイ - フレーズの順序付きシーケンスが、矢印の左側に示されている。そして、デコーディングは、G T S Pとして、グラフのノードがすべての可能ペア (w, b) を表す態様で公式化される。ここで、wはソース文s内のソース単語であり、bはこのソース単語を含むバイ - フレーズである。ここでは、同じ単語タイプでも出現が異なれば異なる単語と考える。特別なバイ - フレーズ $b_{\$} = (\$, \$')$ が導入され、ここで、 $\$$ (または $\$'$) はソース (または目標) 文の開始を示す特別なソース単語となり、かつ、ペア ($\$, b_{\$}$) に関連付けられる、対応する追加のグラフ・ノード $\$ \$ = (\$, (\$, \$'))$ が導入される。

30

【0034】

グラフ・クラスタは、共通のソース単語wを共有するグラフ・ノードのサブセットになり、ノード $\$ \$$ は、ソース単語 $\$$ に関連付けされたクラスタ内の唯一のノードになる。グラフのノードMとNとの間における遷移のコストは、次のように定義される。Mが (w, b) の形式であり、Nが (w', b') の形式であり、bが単一のバイ - フレーズで、wおよびw'がb内で連続する単語である場合、遷移コストは0である (遷移コストがない)。直観的に言えば、bの最初の単語の使用に一旦掛かり合えば、bによってカバーされるほかのソース単語に移動するための追加のコストはない。Mが (w, b) の形式であり、wがバイ - フレーズb内の「一番右のソース単語」であり、Nが (w', b') の形式であり、w' = wがb'内の「一番左のソース単語」である場合、遷移コストは、バイ - フレーズbを選択した直後にバイ - フレーズb'を選択する実際のコストに対応する。ソース文から見ると、これはbのソース側を消費した後のb'のソース側の「消費」に対応し (ソース文内のそれらの相対的なポジションによらない)、目標側から見ると、これはbの目標側を生成した直後におけるb'の目標側の生成に対応する。

40

【0035】

遷移コストは、その場合、バイ - フレーズ・ライブラリ 20 (図1) 内のbに関連付け

50

された静的コストを含むいくつかの寄与の和になる。このコストは、順方向および逆方向の条件付き確率、パイ・フーズ内の目標単語の数、およびこれらの類（導入部分の説明を参照されたい）といった構成要素に対応する。「歪み」コストは、ソース単語 w を消費した直後にソース単語 w' を消費する選択に関連付けされる。 w' がソース文内において w に直接続く単語である場合には、このコストがゼロであり、 b および b' の目標側の連続性がそれらの目標によって与えられる状況に対応する。そのほかの場合には、ソース文内の w および w' のポジションを $pos(w)$ および $pos(w')$ とするとき、 $(pos(w) + 1 - pos(w'))$ の絶対値としてコストが計算される。「言語モデル」のコストは、 b の目標単語をもたらしたばかりの文脈において目標単語 b' をもたらすコストである。バイグラム言語モデルが仮定される場合には、 b および b' がわかるとすぐにこのコストをあらかじめ計算することが可能になる。というのも、これは b がその目標側に少なくとも1つの単語を含み、そのことが b の最後の目標単語を知った上での b' の最初の目標単語の寄与の計算を可能にするからである。 b' の2番目、3番目等々の目標単語については、 b' のみを基礎として寄与が計算される。注意されたいが、このバイグラム・モデルの制限は、ここで論じられているほかの手法を使用して克服できる。

【0036】

パイ・フーズ b および b' のうちの1つが $$$$$ に等しい場合には、以前の寄与の簡単な適応を容易に実行することができる。ほかのすべての場合には遷移コストが無限となり、換言すれば M と N の間のグラフ内にエッジが存在しない。

【0037】

図5は、既存のエッジのサブセット、すなわちノード `translation-mt` に入るか、または出るすべてのエッジ82だけが示された、ソース文「`cette traduction automatique est curieuse`」についての遷移グラフ80を図解している。注意を要するが、`translation-mt` の唯一の後続ノードは `automatique-mt` であり、`cette-ht` は、`translation-mt` の先行ノードでない。それに代えて、`automatique-m` および `automatique-a` から `translation-mt` にエッジを引くことは可能であるが、その種のエッジは、実際のところ、`translation-mt` からの唯一の出口が `automatique-mt` であり、このノードがそのクラス内のほかのノードと排他であることから、横切ることができない。

【0038】

図6Aおよび図6Bは、示された2つの出力に対応する2つのGTS P巡回を図解している。図6Aにおいては、巡回90が、結果として出力 `h.mt.i.s` をもたらす。図6Bにおいては、巡回92が、結果として出力 `ht.s.i.a` をもたらす。

【0039】

上述の図面に関して述べたモデルは、一般TSPの非対称バージョンに対応する。この再公式化を前提とすると、追従可能ないくつかのストラテジが存在し、GTS P用に特別に設計されたアルゴリズムを使用してもよいし、GTS PをATSPに変換し、ATSP用に設計されたアルゴリズムを使用してもよいし、かつ/またはATSPをSTSPに変換し、STSP用に設計されたアルゴリズムを使用してもよい。各オプションは、それぞれ独自の利点および欠点を有する。コンコード(Concorde)ソルバ(たとえば、`www.tsp.gatech.edu/concorde` 参照)等の既存の効率的なTSP用のソルバが使用される場合には、STSP公式化が採用される。しかしながらATSPがSTSPに変換される場合には、TSPグラフ内の頂点の数が2倍になる。さらにまたGTS PからATSPへの経路は、より一般的な公式化が採用されることから非能率の潜在的原因である。したがって、コンコード・テクニクとともにSTSP再公式化を使用することが望ましい。

【0040】

別の重要な要因は、厳密な解が望ましいか、または近似解で充分とし得るか、ということである。STSPの場合においては、厳密な解法(たとえば、コンコード・ソルバ)を

10

20

30

40

50

採用すること、または近似アルゴリズム（たとえば、リン・カーニハン（Lin-Kernighan）のヒューリスティック）を使用することができる。

【0041】

言語モデルがバイグラム・タイプである場合、説明してきたモデルは重要な「マルコフの」性質、すなわち経路のコストは、その経路上の2つの連続するノードの間における遷移のコストに関する加法であるという性質を有する。図6Aにおいて、翻訳候補「this.machine.translation.is.strange」のコストは、単語「is」に関する単語「strange」の条件付き確率を考慮に入れるだけでよく、単語「translation」および「is」に関しては考慮しなくてよい。

【0042】

別の実施態様においては、モデルの性能が、バイグラム言語モデルの使用から、3-グラム言語モデル等のより強力なn-グラム言語モデルの使用に拡張され、いくつかのアプローチを適用することができる。第1のアプローチは、少なくとも2つの単語の目標側を有するパイ・フレーズだけを保持するために、目標側が1つの単語だけを含むすべてのパイ・フレーズを「編集により除外（compiling out）」することを包含する。この態様においては、2つのパイ・フレーズbおよびb'の目標側が連結されるとき、bが少なくとも2つの単語を含むことから、トリグラム言語モデルが、bに関するb'の寄与の計算に十分な文脈を有する。適正な機能を保証するためにパイ・フレーズの概念の拡張を採用し、ここでパイ・フレーズの順序付きシーケンス

【数2】

$$[(\tilde{s}_1, \tilde{t}_1), (\tilde{s}_2, \tilde{t}_2), \dots, (\tilde{s}_k, \tilde{t}_k)]$$

として拡張パイ・フレーズを定義する。ここで、k ≥ 1であり、各

【数3】

$$\tilde{s}_i$$

または、各

【数4】

$$\tilde{t}_i$$

は、ソース単語または目標単語のリストである。k = 1の場合には、この手法はパイ・フレーズのオリジナルの概念に戻る。ソース文sの翻訳のための概念の解釈は、オリジナルの場合といくらか異なる。つまり、それぞれの個別の

【数5】

$$\tilde{s}_i$$

内のトークンは、s内において連続的にマッチングされる必要があるけれども、

【数6】

$$\tilde{s}_1, \tilde{s}_2, \dots, \tilde{s}_k$$

が連続的にマッチングされることは必要ないし、あるいは、sの内側におけるその順序でのマッチングさえ必要ない。これに対して目標側においては、

【数 7】

$$\tilde{t}_1, \tilde{t}_2, \dots, \tilde{t}_k$$

内のトークンは、連続的に、かつその順序で作られる。この概念の下に、前述と同じ可能バイ - フレーズの表を使用して、次に示す拡張バイ - フレーズ mti 、 ti 、および si が提供される。

【0043】

$mti = [mt.i] = [(traduction\ automatique, machine\ translation).(est, is)]$

10

$ti = [t.i] = [(traduction, translation).(est, is)]$

$si = [s.i] = [(curieuse, strange).(est, is)]$

【0044】

これらを使用して、次に示す翻訳を生成することができる。

$[h].[mt.i].[s]$ this machine translation is strange

$[h].[c].[t.i].[a]$ this curious translation is automatic

$[ht].[s.i].[a]$ this translation strange is automatic

20

【0045】

翻訳プロセスのオリジナルのアカウントと現在のそれの間の主要な相違は、拡張バイ - フレーズが、ユニットのシーケンスを通じて以前に達成できたものを単一のユニットの下にカプセル化することである。GTS P グラフとしての翻訳プロセスのエンコードについては、それが簡単になり、グラフのノードがこの場合はペア (w, b) であり、それにおいて w はソース文の単語、 b は拡張バイ - フレーズ、クラスは同一の w を有するノードのサブセットである。オリジナルの規則に対する拡張によって、 w が

【数 8】

30

$$\tilde{s}_1$$

の最初の単語であり、形式

【数 9】

$$(w, b) = (w, [(\tilde{s}_1, \tilde{t}_1), (\tilde{s}_2, \tilde{t}_2), \dots, (\tilde{s}_k, \tilde{t}_k)])$$

であるノードに入るグラフ内の経路は、それが拡張バイ - フレーズ b を「離れる」(この時点において実際の選択、すなわち次の拡張バイ - フレーズの選択がある)前に、

【数 10】

40

$$\tilde{s}_1$$

のすべての単語を順番に通じ、続いて

【数 11】

$$\tilde{s}_2$$

のすべての単語を通じ、同様に繰り返して最後に

【数 1 2】

$$\tilde{s}_k$$

のすべての単語を通らなければならない。拡張バイ - フレーズ

【数 1 3】

$$[(\tilde{s}_1, \tilde{t}_1), (\tilde{s}_2, \tilde{t}_2), \dots, (\tilde{s}_k, \tilde{t}_k)]$$

の「内側の」コストは、

10

【数 1 4】

$$(\tilde{s}_1, \tilde{t}_1)$$

から

【数 1 5】

$$(\tilde{s}_2, \tilde{t}_2)$$

等のように遷移するときに生じるであろうコスト（歪みコストを含む）を加算することによって事前編集できる。概して言えば、基本バイ - フレーズの構成要素にわたって対応する経路を考慮することにより生じるコストを「回収する」ことによって拡張バイ - フレーズにわたる経路のコストを計算することは簡単である。

20

【0 0 4 6】

バイグラム言語モデルからトリグラム言語モデルへの移動の問題に戻るが、バイ - フレーズ・ライブラリから単一単語の目標を有するバイ - フレーズ i を取り除くステップ、および拡張バイ - フレーズ mt_i 、 ti 、 si 等々（たとえば、ライブラリ内のバイ - フレーズの i との連結からなるすべての拡張バイ - フレーズ）をライブラリに追加するステップが実行される。これらの拡張バイ - フレーズは、すぐ次にもたらされる目標単語（与えられている例においては、それぞれ単語「strange」、「automatic」、および「automatic」）についてのトリグラム確率を計算する十分な文脈を提供する。これらのステップが、 i と類似の、すなわち単一単語の目標を有する（手元のソース文に適切な）すべてのバイ - フレーズについて網羅的に実行されると、各ポイントにおいてトリグラム言語モデルの計算を可能にする表現が得られる。

30

【0 0 4 7】

図 7 は、バイ - フレーズ「 i 」だけが取り除かれた、現在は「 i 」をカプセル化している拡張バイ - フレーズのいくつかだけがグラフを通る 1 つの有効な回路または巡回を定義するエッジ 102 を伴って示されているグラフ 100 を図解している。気付かれるであろうが、 mt_i が十分に大きな目標文脈を提供することから、2 つのノード（ est 、 mt_i ）および（ $curieuse$ 、 s ）を接続するエッジが、ここでトリグラムのコスト p （ $strange | translation \ is$ ）に関連付けされる。

40

【0 0 4 8】

図 8 は、選択的なオンデマンド絞り込みを伴う第 2 のアプローチを図解している。今述べたばかりの網羅的な「編集により除外」する方法は、基本的に有効であるが、翻訳されるべき文について m 個の関連するバイ - フレーズがあり、そのうちの k 個が単一単語の目標を有している場合には、 $k \cdot m$ 個の拡張バイ - フレーズが作られ、 k が m に関して大きくなれば直ちに TSP ソルバについての有意のオーバーヘッドを表すおそれがある。この効果は、編集により除外する方法が $n > 3$ を伴う n グラム言語モデルに拡張されると悪化することがある。

【0 0 4 9】

50

この効果を緩和するために、第2のアプローチは、2つの構成要素を有する選択的絞り込みを使用する。第1の構成要素は、何らかの広い評価基準（長さ1の目標側を有する等）に関してすべてのノードの文脈を絞り込むのではなく、グラフ内の選択されたノードの文脈を絞り込む能力である。その種の絞り込みは、グラフ内のノードの少数派のためのトリグラム文脈を提供するが、残りのノードについてはバイグラム文脈だけとなる。第2の構成要素は、その種の絞り込みのための最適TSP解と、グラフ内のすべてのノードについてトリグラム文脈が使用された場合に到達する「真の」最適解との間の結合不等式の維持からなり、真の最適解への絞り込みプロセスの収斂を保証する。

【0050】

したがって図8に、G T S Pグラフ110の選択的絞り込みを図解する。G T S Pグラフ110において、a、b、およびcはグラフ内の、異なるクラスタに属するノードであり、ノードbに出入りするエッジ（ノードを接続するライン）のうちのいくつかが表示されている。それに加えて、各エッジについてのコストまたは重みが、 c_1 、 c_2 、 c_3 、 c_4 、および c_5 とラベル付けされて示されている。変換後のG T S Pグラフ112においては、ノードbが、bと同じクラスタに属する（したがって、相互に排他的な）2つの「クローン」ノードb1およびb2に置き換えられ、bのすべての点に関してまったく等しいが、異なる入射エッジを有し、aに関する入射エッジ（そのうちの1つだけが示されている）は変化しないが、新しくcに到来するエッジが追加されている。ここではノードb1を「aの直接の後続ノードであるという文脈におけるb」として解釈することが可能であり、b2は「そのほかの任意の文脈におけるb」として解釈される。

【0051】

この変換において注意する最初の性質は、コスト c_1 および c_2 が c_3 に等しいと仮定された場合に、変換後のグラフ112が1つのノードおよび3つのエッジをオリジナルより多く有するが、注意深い観察によってわかるとおり、最適巡回は正確に同じであり、同一のコストを伴うことである。しかしながらここでは、b1（またはb2）がaの直接の後続ノードである（または、直接の後続ノードでない）という文脈に特化され、したがってこれらの特化された文脈が、この追加の知識に関してコスト c_1 および c_2 をより良好に定義するために利用される。特に、この変換がSMTの状況に適用されるとき、 c_1 はaに関連付けされた目標単語を承知しており、cの最初の目標単語を条件付けするためにトリグラム言語モデルを利用することを可能にする。

【0052】

G T S Pグラフを前提として、グラフに関する良好に形成された巡回が考慮される。G T S Pグラフのエッジに対して与えられる重みに従って、 c_1 が特定のコスト、すなわち見かけのコストを有する。何らかの外部測定に従って、同一の巡回が異なるコスト、すなわち真のコストを実際に有することがある。この状況の例は、G T S Pエッジが言語モデルのためのバイグラム・コストを持つが、真のスコアはトリグラムの知識に従って計算されるべきであるときに生じる。より一般的に言えば、巡回の真のコストは、グラフのエッジに局所的な重みによって計算可能であるより、いくぶん局所的でない巡回の性質に依存し得る。

【0053】

G T S Pグラフのエッジに関するコストは、グラフに関して任意の良好に形成された巡回について、巡回の見かけのコストが真のコストより小さいか、またはそれに等しい場合に限って「楽観的（optimistic）」と定義することができる。「楽観（optimism）」の概念は、ツリー・サーチにおける許容可能なヒューリスティック（たとえば A^* ）の概念と何らかの類似性を有し、その場合においては「現実性のある」楽観的エッジのコストが注目される。標準サーチ・ヒューリスティクスは、サーチ・ツリーの拡張における局所的な決定を得るために使用されるが、ここで述べているヒューリスティック手順は、問題グラフのより一層正確な明細の反復的な提供に焦点を当てつつ、TSPソルバの「注意」が焦点されるグラフの部分を強調し、続いて「グローバル」TSPソルバに現在の最良の見かけの解を見つけさせる。この概念を踏まえて、実行される一般的な

手順を、次に図 9 A および図 9 B に関して説明する。

【 0 0 5 4 】

図 9 A は、巡回セールスマン問題としてフレーズ - ベースの統計的機械翻訳を実行するための方法を図解している。120 において SMT が GTS P として定義され、GTS P グラフが生成される。122 においては、ソース文に整合するバイ - フレーズが検索され、検索されたバイ - フレーズは、それぞれ GTS P グラフ内のノードに対応する。各バイ - フレーズは、第 1 の言語のフレーズ（たとえば、入力された文の言語において、入力またはソース文内のフレーズに整合するフレーズ）および第 2 の言語のフレーズ（たとえば、第 2 のまたは目標言語に翻訳された入力フレーズ）を含む。124 においては GTS P が解かれる。バイ - フレーズが、GTS P の解の関数として選択され、126 において、
10
選択されたバイ - フレーズの目標または第 2 の言語のフレーズが、GTS P の解によって定義された順序で出力される。

【 0 0 5 5 】

1 つの実施態様においては、GTS P を解くことは、GTS P を ATSP に変換すること、ATSP を標準 TSP に変換すること、および TSP を解決して入力文のブロックを翻訳することを含む。TSP の解決は、コンコード・ソルバまたはリン・カーニハンのヒューリスティックまたはこれらの類を使用して実行される。

【 0 0 5 6 】

図 9 B は、巡回セールスマン問題としてフレーズ - ベースの統計的機械翻訳を実行するための代替または追加の方法を図解している。130 において、その巡回の真のコスト
20
に関して楽観的な GTS P グラフ G_0 の初期仕様が、 $i = 0$ となるように初期化される。グラフ内のノードは、セグメント化された文のブロックに関連付けられたバイ - フレーズを定義する。132 においては、GTS P ソルバ・アプリケーション（たとえば、図 1 の TSP ソルバ 23）が起動され、このグラフに関する最適巡回 i （または、近似ソルバが使用される場合にはその種の最適巡回の近似）が獲得される。134 においては、 i の真のコストが（たとえば、 i のすべてのエッジが既知であることから）計算される。 G_i が楽観的であることから、 G_i に関する真のコストは、 i の見かけのコスト C_a より大きくなる。136 においては、見かけのコストと真のコストの間の差 D があらかじめ定義済みの閾値より小さいか否かに関しての判定が行われる。

【 0 0 5 7 】

これら 2 つのコストの間の差 D が特定の閾値より小さい場合には解 i が出力され、
30
138 においてこの方法が終了する。そうでなければ 140 において G_i の少なくとも 1 つのノード、特に、 i 上に（可能性としては、ほかのいくつかにも）現れている特定のノードが、図 7 の原理に従って絞り込まれる。この種の絞り込みの間はグラフ G_i が楽観的にとどまるが、より厳しい値が、 i によって提供された i_1 および i_2 のために提供される。すなわち、 $i_1 >$ 等の制約が提供され、かつ可能性として $i_2 >$ 等の制約も提供される。142 においては、絞り込みの結果として新しいグラフ G_{i+1} が獲得される。方法は 134 に戻るが、 $i := i + 1$ を伴う。

【 0 0 5 8 】

この図 9 B の方法は、いくつかの重要な性質を有する。たとえば、任意の反復において
40
、 i の見かけのコストは、元のグラフにおける真の最適巡回 $_{true}$ の新しいコストの下限になる。真の最適巡回 $_{true}$ は、すなわち、すべての巡回の真のコストにわたって真のコストが最小となる巡回である。厳密 TSP ソルバの場合においては、 $_{true_cost}()$ $_{true_cost}()$ $_{true}$ であり、すべての巡回 i について（ $_{true}$ の定義により）、 $_{true_cost}()$ $_{true}$ $_{apparent_cost}()$ i である。なぜなら、 i は、すべての巡回のコストの楽観的仕様より最適であり、特に $_{true_cost}()$ $_{true}$ $_{apparent_cost}()$ $_{true}$ $_{apparent_cost}()$ i であるからである。

【 0 0 5 9 】

G_i および i においてアルゴリズムが終了するときは、 $_{apparent_cost}$
50

$(i) + \text{true_cost}(i)$ である。したがって、 $\text{apparent_cost}(i) + \text{true_cost}(i) - \text{true_cost}(\text{true_apparent_cost}(i))$ である。言い替えると、反復の間に見つめられた巡回 i は、真の最適巡回のそれと無視できる程度に異なる真のコストを伴った真の最適巡回の近似である。

【0060】

アルゴリズムの終了特性については、有限数の絞り込みしか存在しないためにグラフを無限に絞り込むことが不可能であること、および与えられた巡回におけるノードが十分に絞り込まれるとき、巡回の見かけのコスト (c_1 および c_2 等の絞り込まれた重みに依存する) がその真のコストに等しくなることという2つの要因に依存する。

10

【0061】

図9Aおよび図9Bに図解されている方法は、ここに述べられているほかの手法またはアルゴリズムに加えて、コンピュータ上において実行できるコンピュータ・プログラム製品として実装できる。コンピュータ・プログラム製品は、ディスク、ハード・ドライブ、またはこれらの類といった、コントロール・プログラムが記録されたコンピュータ可読記録媒体(たとえばメモリ13)とすることができる。一般的な形式のコンピュータ可読媒体は、たとえば、フロッピー(登録商標)ディスク、フレキシブル・ディスク、ハードディスク、磁気テープまたは任意のそのほかの磁気記憶媒体、CD-ROM、DVDまたは任意のそのほかの光学媒体、RAM、PROM、EPROM、FLASH-EPROMまたはそのほかのメモリ・チップまたはカートリッジまたは任意のそのほかの、コンピュータによる読み出しおよび使用が可能な有体の媒体を含む。それに代えてこの方法を、たとえば無線波および赤外線データ通信の間に生成されるような音響または光の波等の送信媒体およびこれらの類を使用するデータ信号としてコントロール・プログラムが埋め込まれる送信可能な搬送波内において実装することができる。

20

【0062】

例示的な方法は、1つまたは複数の汎用コンピュータ、専用コンピュータ(1つまたは複数)、プログラムされたマイクロプロセッサまたはマイクロコントローラおよび周辺集積回路要素、ASICまたはそのほかの集積回路、デジタル信号プロセッサ、ハードワイヤード電子または論理回路、たとえばディスクリート素子回路、PLD、PLA、FPGA、グラフィック・カードCPU(GPU)、またはPALといったプログラマブル論理デバイス、またはこれらの類の上において(たとえば、プロセッサ12により)実装できる。概して言えば、有限状態マシンの実装が可能であり、続いてそれが図9Aおよび図9Bに示されたフローチャートを実装できる任意のデバイスを、GTSモデルを使用するフレーズ・ベースのSMTを実行するための方法の実装に使用することが可能である。

30

【0063】

図10は、トリグラム言語モデルを使用してフレーズ・ベースの翻訳を実行するためのグラフ150を図解している。認識されることになるが、同じアプローチを4-グラム、5-グラム等々に簡単に適用することができる。トリグラム言語モデルは、単語 z が2つの単語 x および y に続く確率 $p(z | xy)$ の評価を提供する手順を容易にする。1つの実施態様においては、言語モデルがすべての3成分要素(x, y, z)を、それらの確率とともに内部的に記憶する。別の実施態様においては、言語モデルが、明示的に特定のトリグラム、バイグラム、およびユニグラムのためのコーパス・カウントを記憶し、それらのテーブルから $p(z | xy)$ の計算のためにスムージング手法を頼る。

40

【0064】

このアプローチは次のとおりとなる。 $p(z | xy)$ を考慮して目標言語モデルについてのグラウンド・トゥールズが提供される一方、初期グラフのすべてのエッジにトリグラム・コスト

【数 1 6】

$$\text{tri}(z|xy) \triangleq -\log p(z|xy)$$

がラベル付けされる必要はない。しかし、むしろいくつかのエッジ（ノードが目標側に 1 つの単語だけを有するエッジ（ a, b ））に、次のとおり定義される「バイグラム・プロキシ」 $\text{bi}(z|y)$ をラベル付けすることができる。

【数 1 7】

$$\text{bi}(z|y) \triangleq \min_x -\log p(z|xy)$$

10

言い替えると、プロキシ $\text{bi}(z|y)$ は、 y に先行し得る単語 x に関して最大限に樂觀的な y と z の間における遷移のコストのための評価である。ここで最小（ \min ）は、翻訳されるべき特定の文に関係する x に関して求められる。たとえば、その文を翻訳するためのバイ・フレーズの最後の目標単語と同一の x に関して求められ、語彙内のすべての可能な単語 x に関して求められることはない。注意を要するが、 $p(z|y)$ が $p(z|xy)$ から導かれたバイグラム言語モデルを表すとき、 $\text{bi}(z|y) = -\log p(z|y)$ は概して真にならず、したがって $\text{bi}(z|y)$ は、その用語の通常の意味におけるバイグラム確率を表さない。

【0065】

20

したがって、図 10 のグラフ 150 は、図 5 とまったく同じ GTS P グラフを伴って開始することによって獲得され、それにおいてエッジ上の言語モデルのコストは、それらのエッジ上で得られる文脈を前提として可能な限り特定のであり、言い替えるとそれらは、先行するバイ・フレーズに応じて $\text{tri}(z|xy)$ の形式または $\text{bi}(z|y)$ の形式のいずれかとなる。特定の巡回のためのすべてのエッジが示され、それに加えて説明の焦点が当てられるノード est-i 上の入射エッジが略式に（破線で）示されたグラフ 150 によって一例を示す。

【0066】

TSP ソルバが起動され、特定のバイ・コストだけでなくいくつかのトリ・コストも使用して見かけの最適巡回が獲得される。グラフ 150 内の巡回の真のコストが、すべての真のトリグラム・コスト $\text{tri}(\text{this}|\$ \$)$ 、 $\text{tri}(\text{machine}|\$ \text{this})$ 、 $\text{tri}(\text{translation}|\text{this machine})$ 、 $\text{tri}(\text{is}|\text{machine translation})$ 、 $\text{tri}(\text{strange}|\text{translation is})$ 、および $\text{tri}(\$ |\text{is strange})$ の計算を伴って計算される。真のコストは（より小さい）見かけのコストと比較され、その差がより小さければ、それ以上の動作がとられる必要はない。

30

【0067】

図 11 は、図 7 および図 9 に関して述べた手順を使用し、差がより大きい、または等しい場合に得られる、識別された巡回上のノードの少なくとも 1 つの 3 成分要素（ a, b, c ）に適用されるグラフ 160 を示している。たとえば、図 7 の手順を、 $(a, b, c) = (\text{automatique-mt}, \text{est-i}, \text{curieuse-s})$ とともに適用してグラフ 160 を得ることが可能であり、それにおいて curieuse-s に est-i-1 をリンクしているエッジ上には、そのエッジのトリ・コストの計算に十分な文脈が存在する。注意を要するが、 $\text{bi}(\text{strange}|\text{is})$ は、 est-i-2 を curieuse-s にリンクしている破線のエッジ上に保持されるが、 automatique-mt を除外し、 est-i-2 に先行できるノードにわたって最小化することによるコストの再計算によって、わずかにより厳しい結合を得ることができる。

40

【0068】

グラフ 160 が獲得されると、TSP ソルバが、再起動されて手順が反復的に実行される。新しい見かけの最適巡回がエッジ（ $\text{est-i-1}, \text{curieuse-s}$ ）を含む

50

ときは常に、このエッジのコストが以前のものより正確になる。単に、各反復時に見かけの最適巡回において1つの3成分要素だけを絞り込むことからなるアプローチは、いずれかのポイントで終了する。そうでなければ、いずれかのポイントにおいて、見かけの最適巡回において、必然的に、すべてのエッジがトリ・コストを持つことになり、したがってその見かけのコストが真のコストに等しくなり、そして 閾値の評価基準を満たすことになるためである。

【0069】

これは、アルゴリズムの収斂の形式的な証明を提供するが、変換されたグラフがコストの真の状態をより迅速に「模する」ためには、各反復において、単一の3成分要素より多くを絞り込むほうがより効率的なことがある。認識されるであろうが、その種のアプローチのすべての可能な変形および/または組み合わせは、この説明の範囲内となることが意図されている。しかしながら、単純な方法は、現在の見かけの最適巡回上に現れるすべての3成分要素(a, b, c)を絞り込むことである(ソース文の長さがnであれば、多くともnのその種の絞り込みを行うことが可能である)。上述した編集により除外する手法は、すべてのトリグラムの網羅的な絞り込みに対応し、したがって選択的な絞り込み手法で行うことが可能な1つの極端な場合であることに注意を要する。

10

【0070】

この態様においては、SMTの状況で起こりがちなように、新しい見かけの巡回がいくつかの下位経路を以前の見かけの巡回と共有する場合に、それらの経路に関する絞り込み済みの知識を利用することになる。注意されたいが、たとえ以前の巡回において 閾値条件が満たされていない場合であっても、新しい巡回が実際に以前の巡回とまったく同じであることが可能であり、これは言語モデルのコスト以外のコストが、この経路上におけるバイ・コストからトリ・コストへの移動に関連付けされる補償より大きくなり得ることによる。

20

【0071】

以上の説明はトリグラムの扱いに焦点を当てたが、選択的絞り込みのアプローチがnグラム(nは整数)に拡張できることは容易に理解される。このアプローチは、トリグラム、4-グラム等々のための拡張された状況の提供に採用可能であり、方法は、巡回のいくつかの部分の言語モデルのコストを絞り込むことが、ほかのすべてのSMT制約を前提に巡回の最適性を先細りにしない限り効果的である。

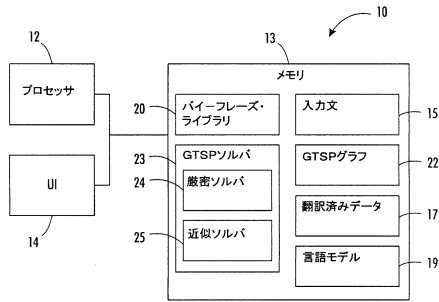
30

【符号の説明】

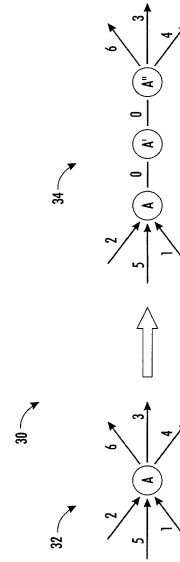
【0072】

10 システム、12 プロセッサ、13 メモリ、14 ユーザ・インターフェース、15 入力文、17 翻訳済みデータ、19 言語モデル、20 バイ・フレーズ・ライブラリ、22 GTS Pグラフ、23 GTS Pソルバ、24 厳密ソルバ・アルゴリズム、25 近似ソルバ・アルゴリズム。

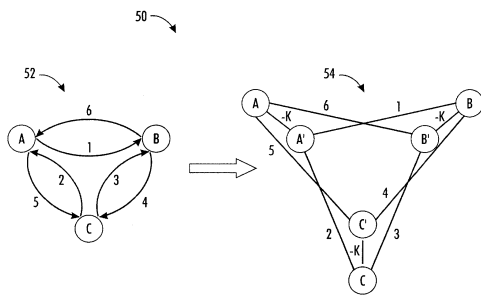
【図 1】



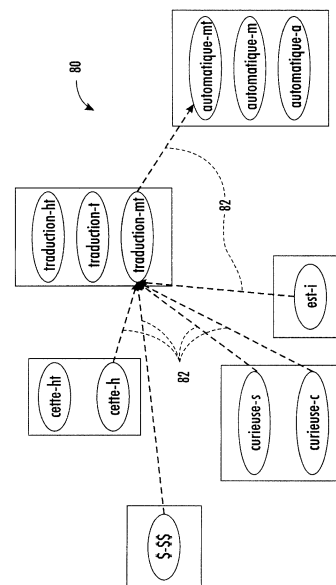
【図 2】



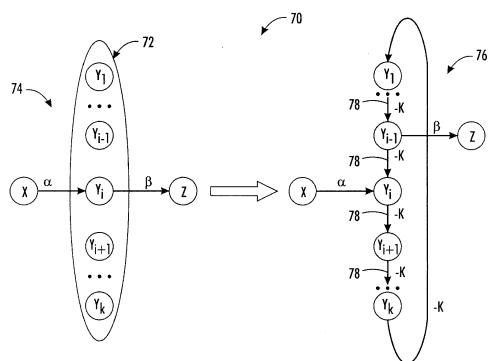
【図 3】



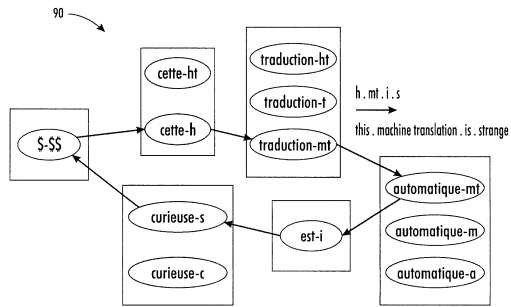
【図 5】



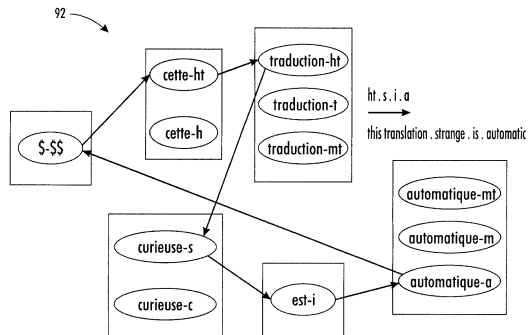
【図 4】



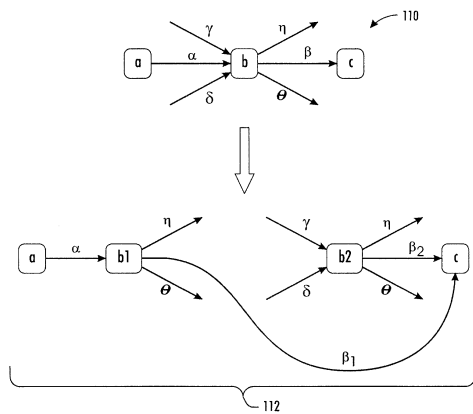
【図 6 A】



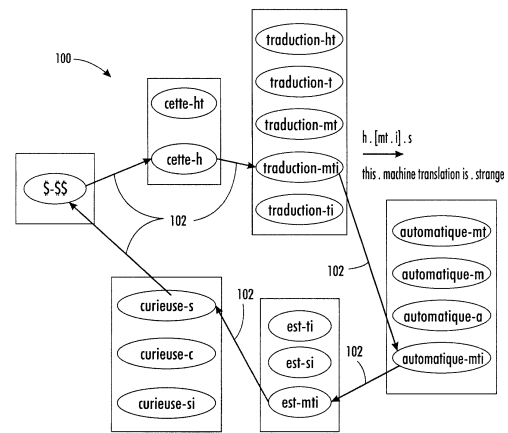
【図 6 B】



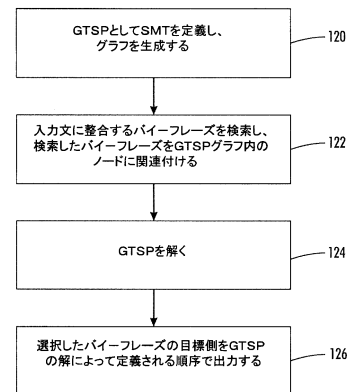
【図 8】



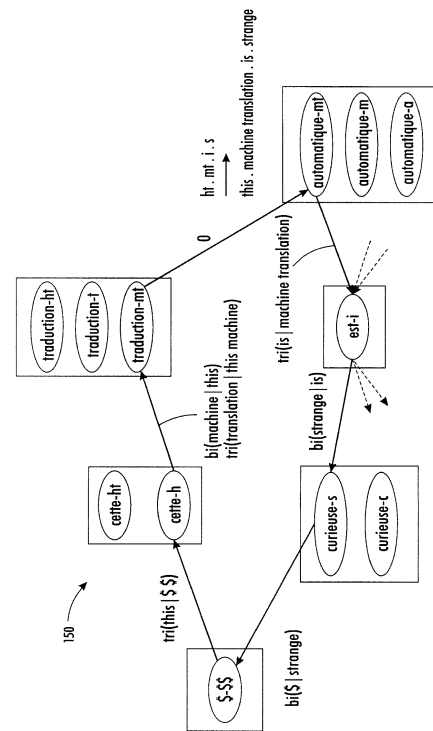
【図 7】



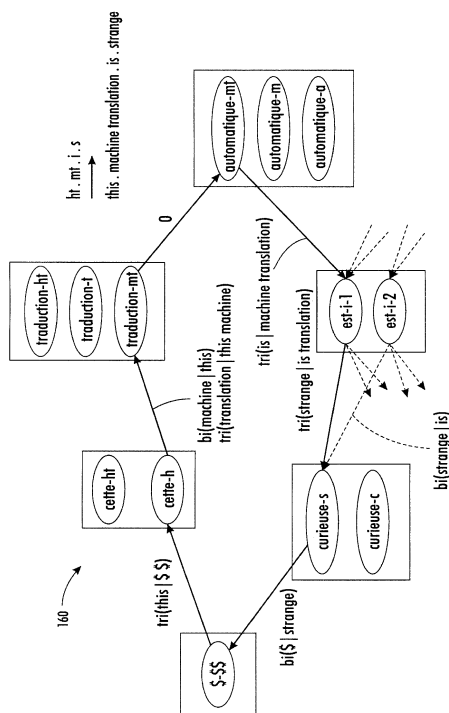
【図 9 A】



【 図 1 0 】



【 図 1 1 】



フロントページの続き

(72)発明者 ニコラ カンチェッタ
フランス グルノーブル クール ジャン ジョレス 31

審査官 長 由紀子

- (56)参考文献 国際公開第02/097663(WO, A1)
特開2007-328483(JP, A)
特開2004-102946(JP, A)
Yookyung Kim et al., Sehda S2MT :Incorporation of Syntax into Statistical Translation System, International Workshop on Spoken Language Translation (IWSLT) 2005, 米国, 2005年 8月25日, p.153-158, URL, http://20.210-193-52.unknown.qala.com.sg/archive/iwslt_05/papers/slt5_153.pdf
本橋 瞬 外2名, Lin - Kernighanアルゴリズムをカオス駆動する巡回セールスマン問題の解法, 電子情報通信学会技術研究報告, 日本, 社団法人電子情報通信学会, 2008年 3月21日, 第107巻 第561号, P.43~48
越川 満 外4名, 統計的機械翻訳におけるフレーズ対応最適化を用いた翻訳候補のリランキング, 言語処理学会第15回年次大会発表論文集, 日本, 言語処理学会, 2009年 3月 2日, p.861-864

- (58)調査した分野(Int.Cl., DB名)
G06F 17/20 - 28