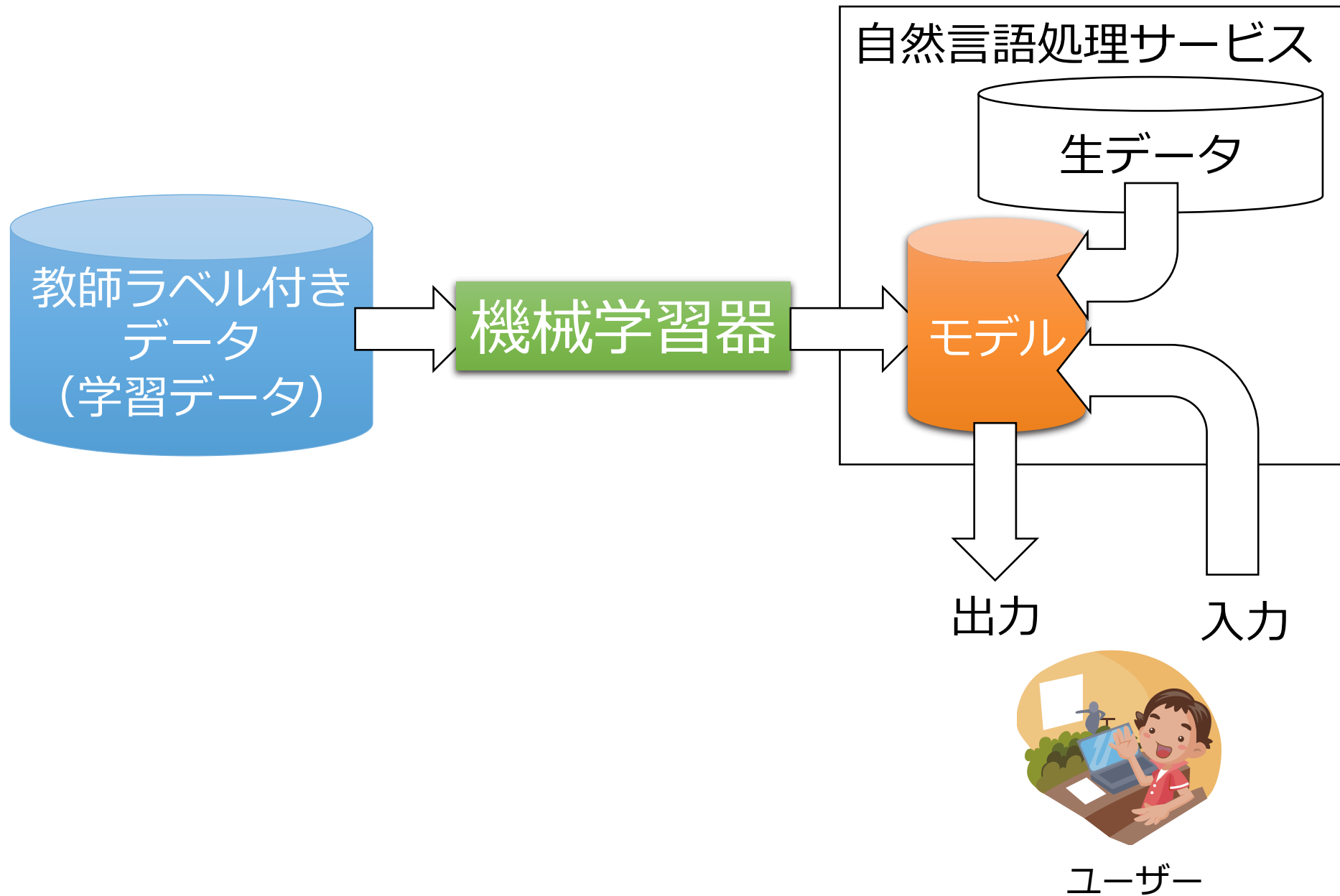


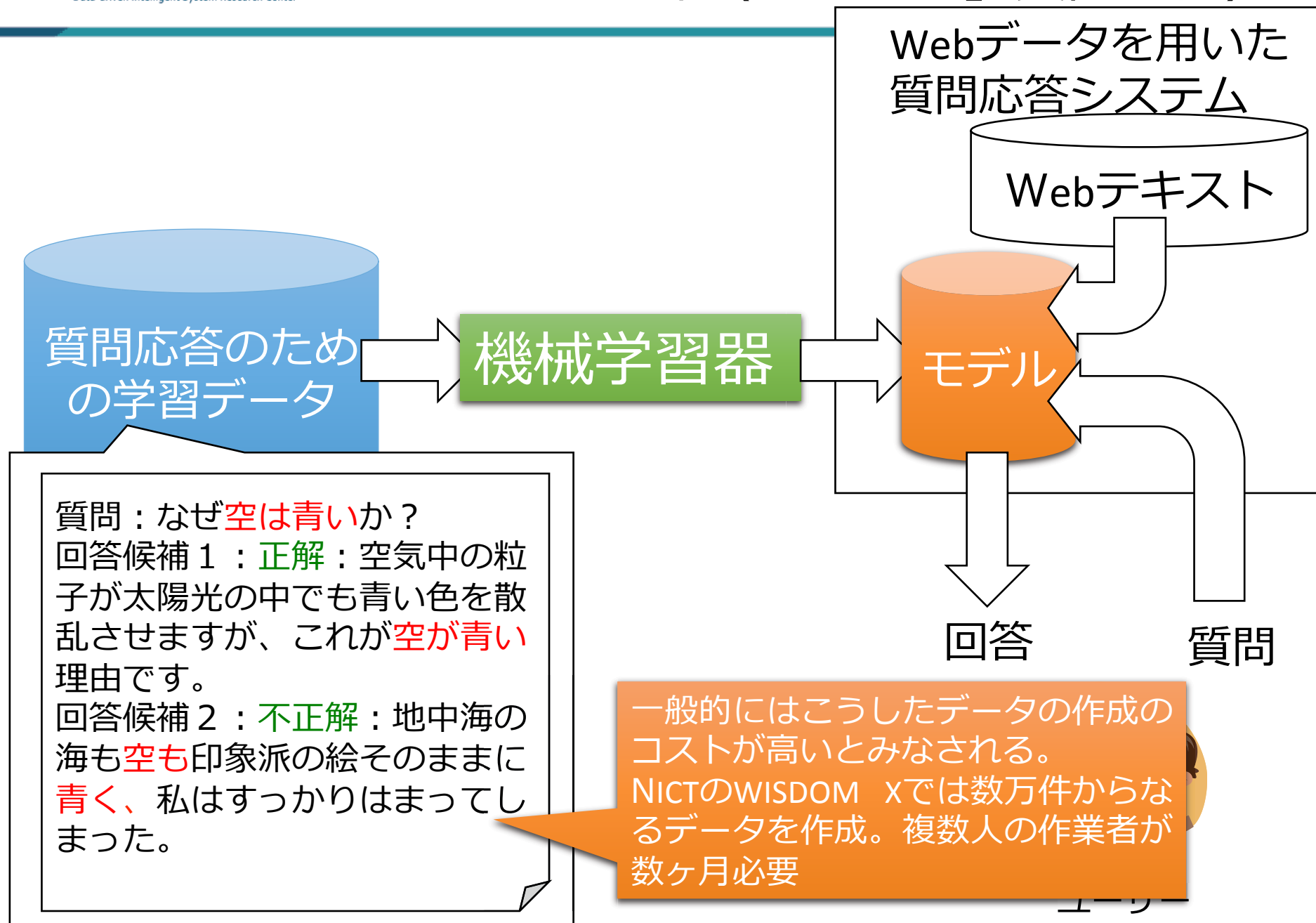
自然言語処理のための データ整備について

鳥澤 健太郎

NICT, DIRECT

2017.2.14





- 二種類のデータ
 - 生データ
 - 人工知能の研究開発での利用を意識した加工が一切なされていないもの
 - 例：生のWebテキスト、科学技術論文
 - 機械学習等でそれだけで有用な場合は「まれ」
 - 例：SNS上のユーザがタグ付けした画像情報
 - 普通は、機械学習で学習したモデルの適用先
 - 例：質問の回答をその中から探してくるテキスト（そのテキストに、学習済みのモデルを適用することで回答を見つける）
 - 教師ラベル付きデータ（いわゆる「学習データ」）
 - 生データに機械学習用のラベルが付与されたもの
 - 機械学習の学習時に学習データとして利用
 - 生データに比して作成には圧倒的にコストがかかる
 - 人がラベルを付与する作業が必要なため
 - 例：質問と、その質問の回答の候補が書かれたテキストのペア
 - 回答候補が正解か、不正解かを示すラベルがふられている

なぜ学習データが大量に必要なか？

- 自然言語では同じことを言う方法が無数にある

回答候補 1 : **正解** : 空気中の粒子が太陽光の中でも青い色を散乱させますが、**これが空が青い理由**です。

回答候補 2 : **正解** : 空気中の粒子が青い色を散乱させるので、**空が青く**なります。

回答候補 3 : **正解** : 空気中の粒子が青い色を散乱させ、**空が青く**なります。

回答候補 4 : **不正解** : 地中海の海も**空も**印象派の絵そのままに**青く**、私はすっかりハマってしまいました。

質問 : なぜ**空は青い**か？

80年代、90年代の対処法
パターン、ルールをプログラマが書く

```
if  
  (これが (質問の内容) 理由です  
  or  
  ...ので、 (質問の内容) 。)  
then 正解  
else 不正解
```

ルール中の「理由」は「原因」で置き換えてもいい。
回答候補 3、4は「理由」のような手がかりがないので、抽出できない....

エンドレスな修正、悩み
⇒開発者が疲れてプロジェクト失敗

• 自然言語では同じことを言う方法が無数にある

回答候補 1 : **正解** : 空気中の粒子が太陽光の中でも青い色を散乱させますが、**これが空が青い理由**です。

回答候補 2 : **正解** : 空気中の粒子が青い色を散乱させるので、**空が青く**なります。

回答候補 3 : **正解** : 空気中の粒子が青い色を散乱させ、**空が青く**なります。

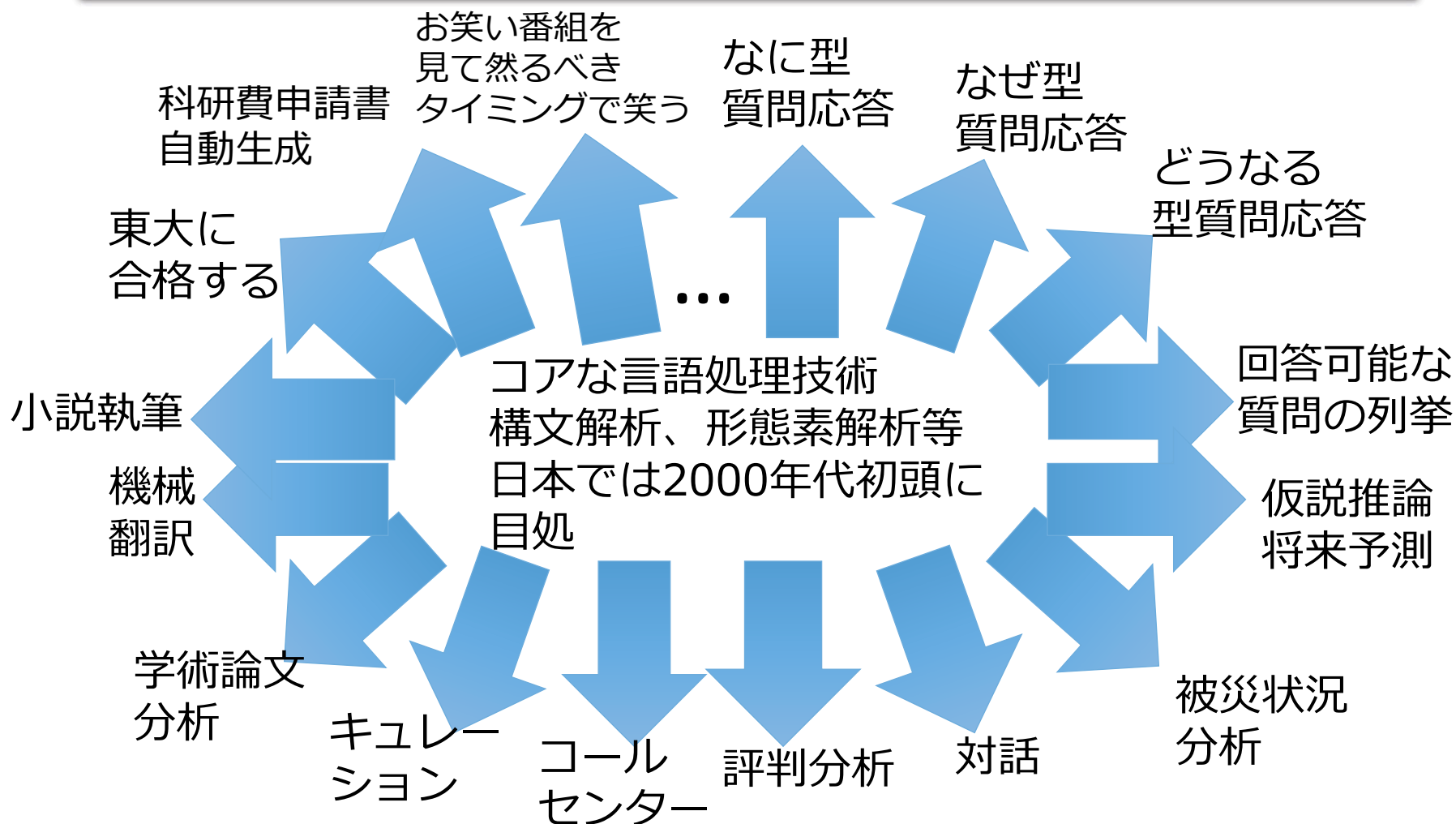
回答候補 4 : **不正解** : 地中海の海も**空も**印象派の絵そのままに**青く**、私はすっかりハマってしまいました。

質問 : なぜ**空は青い**か？

最近の対処法

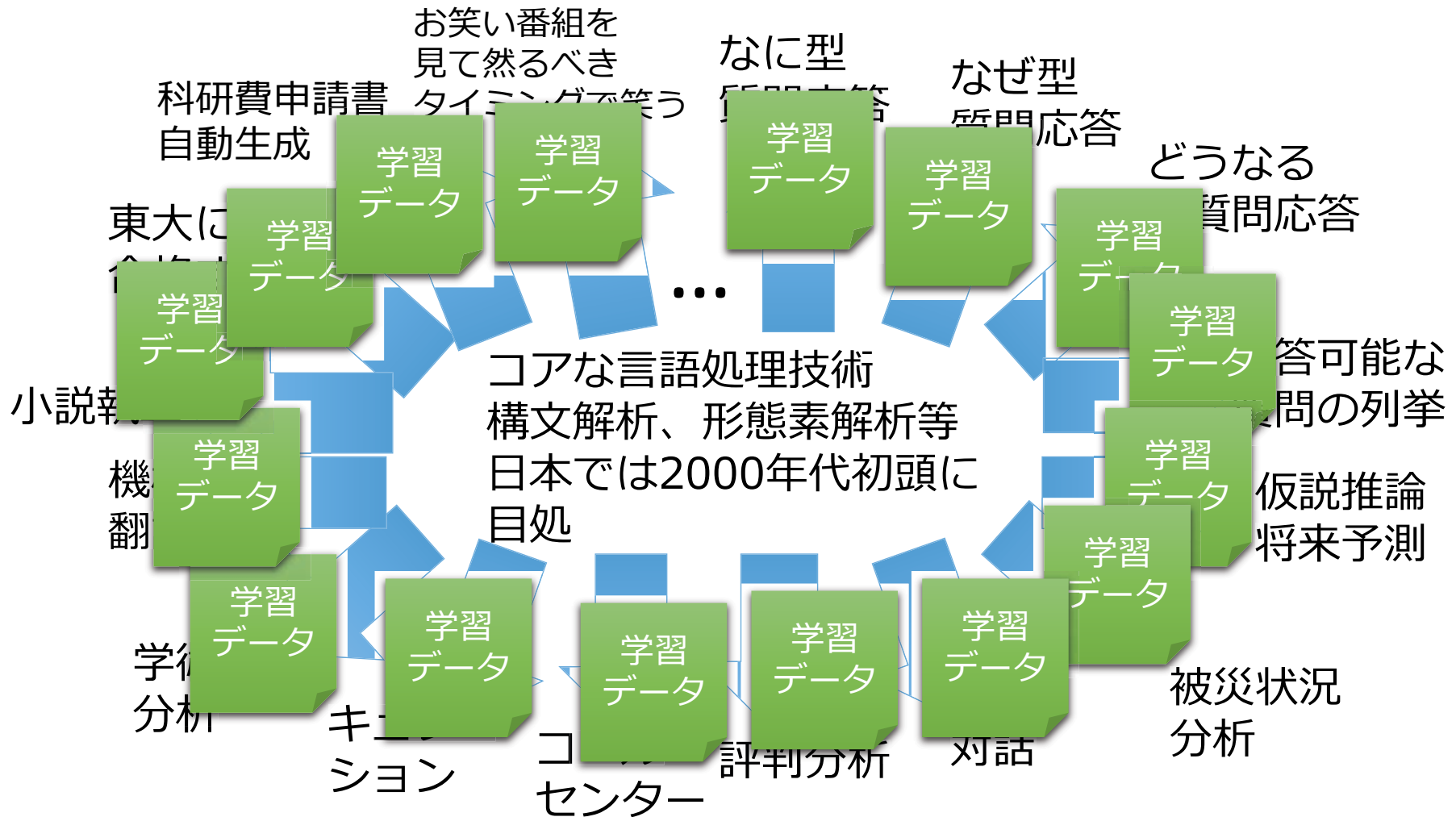
- 疲れるプログラマではなくて、疲れない**機械学習器**に「パターン、ルールのようなもの」(実際には非線形の関数、これを**モデル**と称する)を自動的に書いてもらおう
- ただし、機械学習器には左側にあるような正解と不正解の事例(つまり、学習データ)を**たくさん**与える必要
- **たくさん** = 自然言語の無数のバリエーションをある程度カバーできる量
- パターンを書くのは熟練のプログラマでなければダメだが、正解、不正解の判定は普通の人でも可能⇒コスト的に見合う

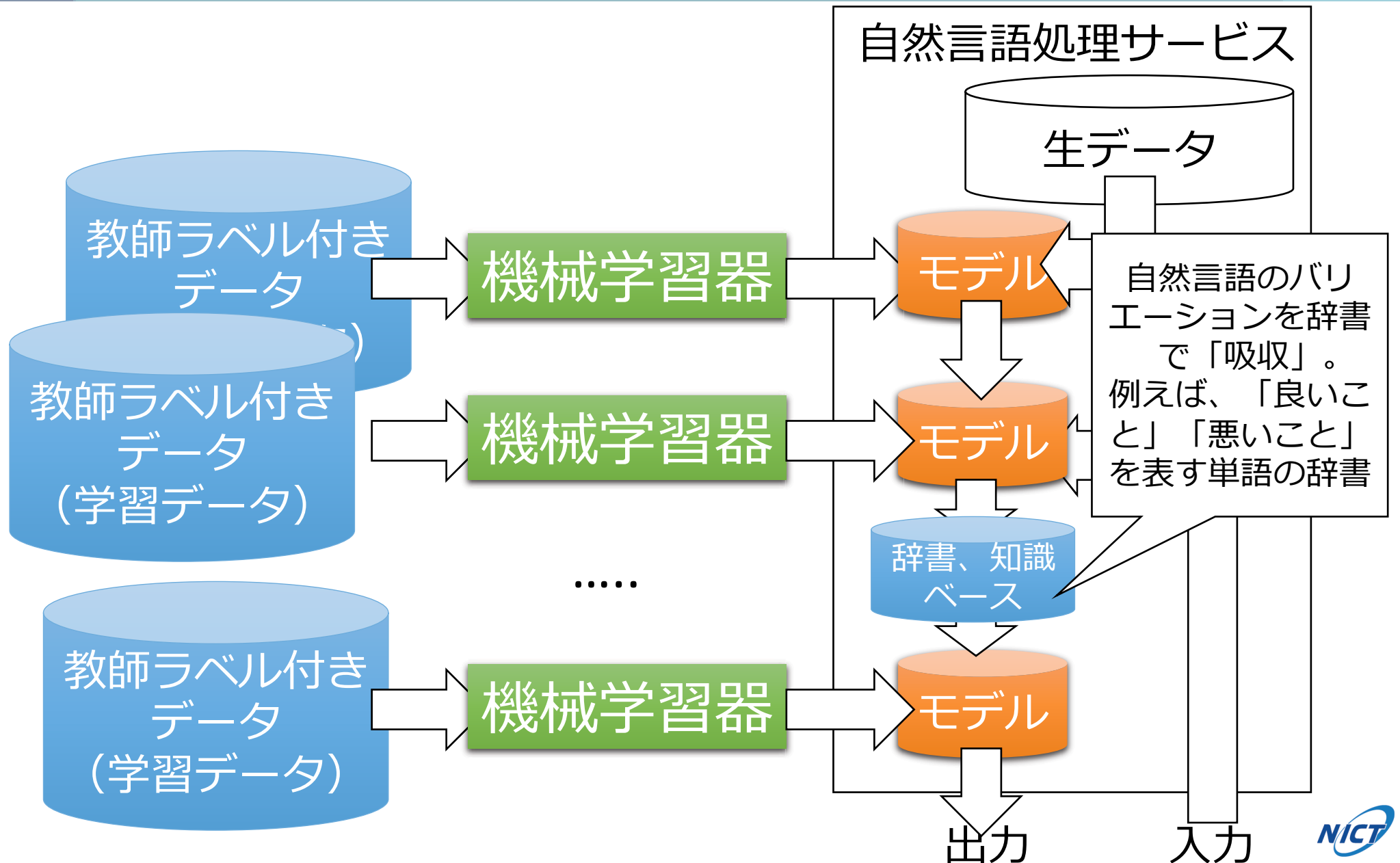
- 日々、新たなタスクが提案され、売り出される
- 現状は、ドラえもんには程遠いが、有望な提案も多数



言語処理研究のタスク

- しかし、当分、タスク毎に学習データは必須





- 学習データの設計は目標とする実サービス、アプリケーションに依存する
 - 分野、タスクの設計に依存
- 学習データの設計は目標とする実サービス、アプリケーションのアーキテクチャにも依存する
 - タスクの分割、速度、前処理等に関する制約にも依存
- この2点に十分注意を払っていないければ、データを作っても、論文が書けるだけで、実運用、社会実装には結びつかない
 - データ作成とサービスのテストはサイクルにならない
- 逆に言うとデータだけがある状況というのは技術の社会実装を阻害する可能性もある
 - 実現不可能なサービスに向けた研究だけが進捗する
 - 一般公開されているデータが論文を書くのにしか使えないのは、実サービスとデータのカップリングが不十分なため
- NICTでは、サービス、アプリケーションの開発を行う研究者が学習データの構築法、設計まで考える
 - サービス、アプリケーションの実装が第一、データは重要だがその手段
 - 場合によっては実ユーザのヒアリング等も参考にしつつ学習データを設計する

- その他に考慮しなければいけないこと
 - そもそも作業者が実行可能なタスクか？
 - 例：医療に関するデータ作成は医療従事者以外に作れるのか？
 - 例：各分野の専門家は十分に確保できるか？
 - 作業マニュアルは十分に分かりやすいか？
 - 例：ゼロ照応の同定を行ってください => 省略されている主語、目的語を補って、意味がとおるものを書いてください
 - 作業者間でラベル付けが一致するか？
 - 作業マニュアルの出来不出来に大きく依存
- データ作成者と自然言語処理研究者、システム開発者が一緒になって上記の点を十分に検討しなければ良いデータは作れない
- ただし、翻訳は翻訳者の基準が確立しているであろうから若干状況は違うであろう

- NICTではWISDOM X、DISAANA等のライセンスを開始
 - 一部データの提供は著作権法の問題から不可能
 - Web等から抽出した長い表現
 - その場合は、データではなく機械学習のモデルや抽出のプログラムを提供
- 学習データ等もサービス、アプリに依存する⇨
ソフトウェアとデータはセットで提供しなければ意味がない
- NICTでは、今後、他組織との連携は人材交流が必須と考える
 - NICTの肩書きがなければ、著作権法等の問題によりNICT内部のデータ(特にWeb由来のもの)を触る事ができない
 - NICT内で作業を行い、成果をライセンス
 - 出向元でもデータ作成は必須。そのノウハウを伝授。
 - どういったタスク、アーキテクチャをターゲットとするかよく話し合いをする必要

- 上記のようなデータ作成者とサービス提供者(データ利用者)がWin-winの関係を構築し、ポジティブなスパイラルを回すにはどうしたらよいか？
 - どのような分野、タスクをアドレスすべきか？
 - サービス提供者は何が実現できると嬉しいか？
 - サービス提供者（データ利用者）はどのようなコミットメントが可能か
 - データ作成側は何が実現できると嬉しいか？
 - 密なコミュニケーション、人材交流、育成？
- ボランティアベース、理想論ベースでは長期的には絶対に上手くいかない
 - データを作る側にも成果（論文？実システム？）が求められる
 - 論文の数を稼ぐだけならば、データを作るという戦略は損
 - 過去の取り組み
- NICTの立ち位置（案）
 - NICT: なるべく幅広い分野に対応するシステム+データ（例：WISDOM X）とパイロットシステムとしての分野特化型システム(例：DISAANA)
 - 他組織：社会的ニーズと直結した分野特化型システム

応用例：車内対話ロボット

対話ロボット：今日の一番のニュースは英国がEUから離脱する…

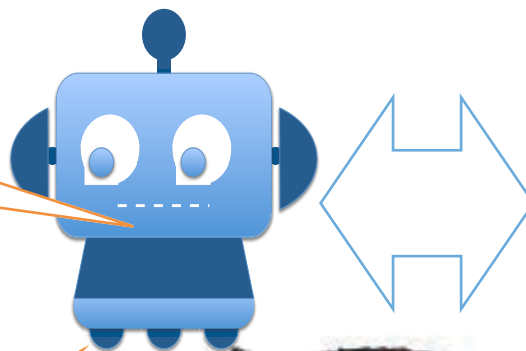
ドライバー：それでどうなるの？

対話ロボット：スコットランドが独立するようです。雇用も失われるとか。英国国教会がカトリックに変わるとかという話まであります。

ドライバー：スコットランドの独立とか、何が起こるか分からないねえ。

対話ロボット：スコットランドが独立すると、ポンド、イギリス株に打撃です。

ドライバー：イギリス株売ります。証券会社に電話して。



通勤中に
一人ブレインストーミング

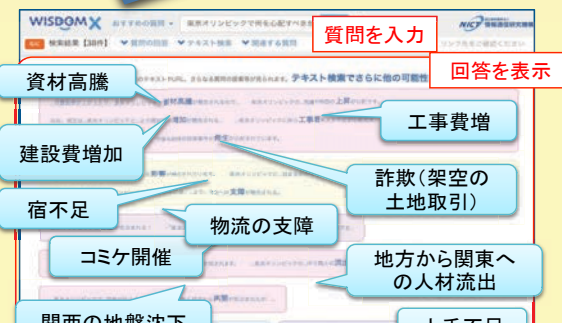
物理的なロボットやセンサーと対話の組み合わせは
NICTの他グループが研究中。ロボカップで複数回優勝

WISDOM X

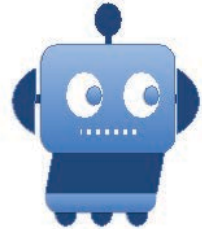
膨大なインターネット上の文字データ（40億ページ、原稿用紙220億枚相当）をもとに、質問に回答や仮説を提供するシステム
（一般公開中）

左の対話例中の情報はすべてWISDOM Xから実際に取得されたもの

質問「東京オリンピックで何を心配すべきか？」



NICT WISDOMくん



iPS細胞で臓器を作るんだって

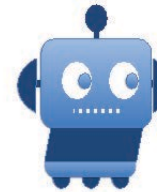
i P S細胞でがんワクチン用細胞量産技術を開発するかも

iPS細胞は様々な可能性を持っているよね

たしかに。i P S細胞でノーベル賞受賞を果たすよね

発話を入力

NICT WISDOMくん



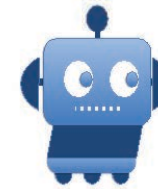
年金支給額が引き下げられるらしいよ

たしかに。気づかれずに年金支給額を減額するんだ

消費者物価が下落気味だからかな

うんうん。給料が増えない

NICT WISDOMくん



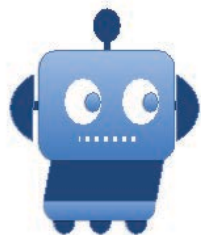
海外旅行に行ってみたいな

なるほど。海外旅行でビーチに行きたいよ

ブーケットのビーチに行きたいな

なるほど。ブーケットのホテルJAL系のツアーでブーケットに行こうよ

NICT WISDOMくん



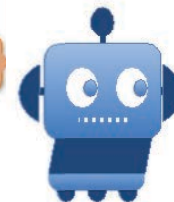
外国人労働者が増えていますね

日本に外国人労働者を入れたい

今、人手不足が深刻ですからね

でも、人手不足は環境の改善に繋がるそうだ

発話を入力



夕ご飯は何にしようかな

きょうの夕ご飯は美味しいイタリアンバーに行くかも

イタリアンは美味しそうだね

イタリアンでジューシーな手づくり燻製を味わうよね

ポイント：深層学習でシステム発話を生成するのではなく、WISDOM X用の質問を生成し、その回答を修正してレスポンスにする

ユーザ入力：ジョブズ、すごいよね

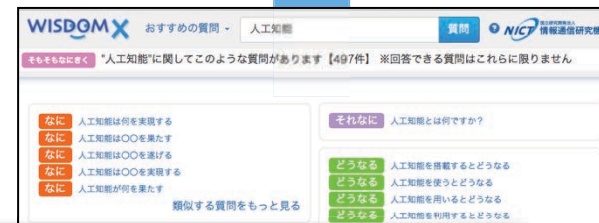
システム出力：Windowsを揶揄するよね～

深層学習のモデルでユーザ入力をWISDOM X用の質問に「翻訳」

回答と元文からシステム出力を合成

質問（複数）
ジョブズが何を見下す？
ジョブズが何を作る？
ジョブズが何に成功する？
....

回答「Windows」
元文「ジョブズがWindowsを揶揄した」
（質問「ジョブズが何を見下す」に対して）



ナイーブな深層学習による対話システムは「うんうん」「あ、そう」といった相槌ばかり言うが、それは防げる

どのような質問を生成するかで、対話システムの「人格」も決まる

例：新聞記事を題材に対話をするとして

質問：南鳥島はどこにある？



システム：南鳥島は**日本最東端**にあります。覚えておきましょう。
(受験生向け対話システム)

新聞記事
「南鳥島沖に球状レアメタル
海洋機構チーム」
日本経済新聞 2016/8/26

システム：自動車会社に影響があるかもしれません。
(ビジネスマン向け対話システム?)

質問：誰がハイブリッド車を製造する



質問：レアメタルは何に使う？



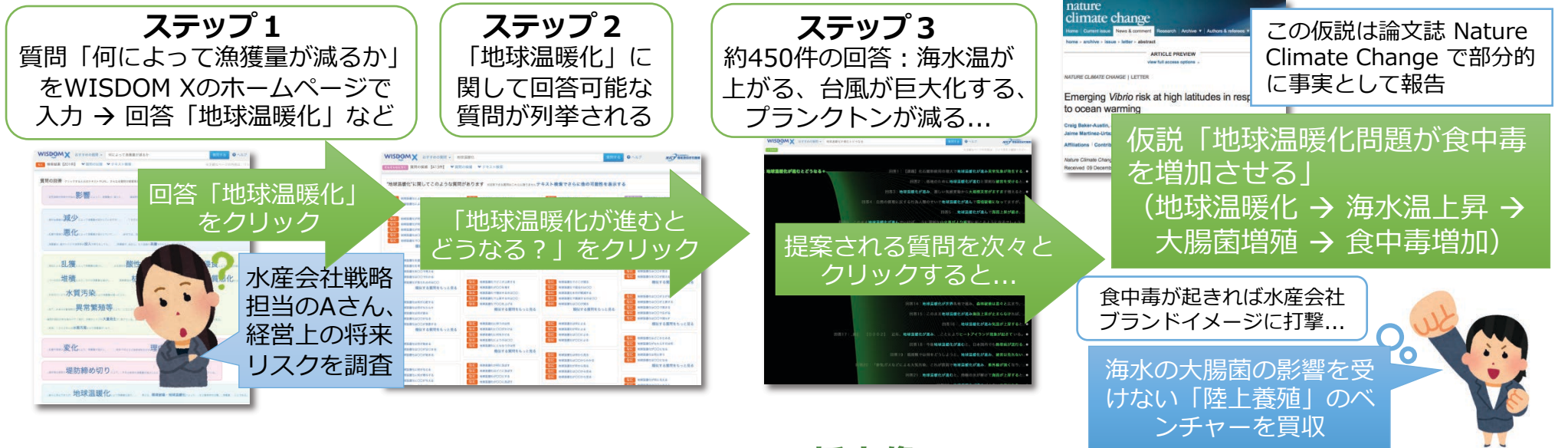
システム：新しい**薩摩焼**とかが現れるかもしれません。
(高齢者の趣味対応用対話システム?)

質問：釉薬を何に使う？



補足：NICTで開発、公開している
システム群

- 自然言語処理技術を用いてユーザーの質問に対してWeb40億ページの情報を基に様々な回答を整理して表示
- また、その回答に対して、システムがさらに質問を追加提案することができるため、ユーザーがその質問と回答をたどることによって、新たな「仮説」を立てることも可能
- 一般公開中



WISDOM Xの将来像

万能対話ロボット(教育、高齢者)



車いすで楽しめるダンスがあるそうです。

ナナフシってオスなしでも繁殖するよ。

民間企業のイノベーション支援



南米でディーゼル油を生成する真菌(水虫の類似物)が発見!

その作戦でいきましょう!

我が社のプラントによく適合しているので、プラントとセットで販売できるかも。

シンクタンク、社会調査



少子化で耕作放棄地が急増!

それでA地方の雇用を増やせますね!

耕作放棄地で行うビジネスには、太陽光発電、魚類の養殖、植物性プランクトンの養殖。A地方に適しているには植物性プランクトンの養殖...

対話技術に応用することで飛躍的に対話機能が向上。ロボットの話に感化されて、ノーベル賞の受賞も夢ではないかも?

民間企業やシンクタンクが活用することで、専門家でなくても、あらゆる技術、出来事、施策の膨大な組み合わせを、人間には実行不可能な規模でシミュレーション可能となり、技術のきっかけにして将来有望な様々なアイデアが生まれる

- SNS（ツイッター）上の災害関連情報をリアルタイムに深く分析・整理して、状況把握・判断を支援し、救援、避難の支援を行う質問応答システム
- ツイートしてから5秒で分析結果を提供可能
- 熊本地震の際には、ツイッター社から人道支援として1ヶ月ほど100%のツイートを提供いただき、その分析結果を提供（平時は10%サンプル）
- WISDOM X, D-SUMMとあわせて民間企業へのライセンスも締結



住民、救援団体からの質問（例：「熊本県で何が不足していますか」）に瞬時に回答

2015年4月より一般公開中

- 熊本地震の際、首相官邸で活用
- 指定避難所以外のニーズ把握
 - 日々変化する要望の把握

↓
熊本県へ指示

平成28年5月11日読売新聞夕刊一面等、報道多数



質問例から選択 ▶ 熊本県で何が不足していますか 検索 関連するツイート中のキーワードを自動的にチェック

検索条件を設定を表示 ←このボタンで絞り込み条件設定を表示できます。

期間指定による検索結果の絞り込み 回答候補の件数：期間内=672件 期間外=37件 最新 2016/04/15 00:00:31 最新 2016/04/19 01:54:57

2016/04/15 20:42:43 自動更新も新しい

回答候補に関する地点を地図に表示中 リックすると回答候補を一覧表示

生活必需品

救援物資

生理用品

回答を地図上に表示し、被災状況を俯瞰可能

回答をピンポイントに抽出


表示解除 衛生用品 (1)

表示解除 ウェットティッシュ (2)

表示解除 毛布 (13)

災害状況要約システムD-SUMM（ディーサム）

- 被災報告の自動抽出結果を整理して要約。リアルタイムで被災状況を把握可能。
- 熊本地震においては、発災後わずか1時間以内には、建物倒壊、負傷者発生、停電、ガス漏れ、信号故障等の被災概況が把握できていた。
- 現地機関からの報告やマスコミ報道を待つまでもなく、初動対応に活用可能。



【災害状況要約レポート（熊本県の被害状況）】2016年6月21日 14:12 自動生成

日時：2016-04-14 (22:25) から 過去 60 分 対象エリア：熊本県

概要：04/14(22:13)熊本市で災害(土砂災害)情報あり。また、04/14(22:22)熊本市でトラブル(水道トラブル)に関する情報、04/14(22:12)熊本市で怪我(負傷)に関する情報を検知しました。

災害: 地震(133), 津波・高潮(18), 土砂災害(1), 火災・火事(1), 水害(1), 風災(1), 悪天候(1), 気象(1), 情報通信機器(1), 怪我: 負傷(7)

対象エリア：熊本県

4月14日21:25 - 22:25のツイート
を要約した画面

※ 熊本地震前震 4月14日(木) 21:26

熊本市 (1726)

益城町 (115)

阿蘇市 (115)

被災報告の多いエリアから表示

情報通信機器（携帯）が繋がらない

ガストラブル（ガス漏れ）が発生

電気トラブル（停電）になる

怪我をする

建造物が崩れる

2016年10月に一般公開開始。開発は内閣府SIPの支援で実施