# Machine Learning System Architecture

Microsoft Translator, a Case Study.

Vishal Chowdhary (@talktovishal, vishalc@microsoft)

# What do we want to do?

# What do you think we do?

# Why is Machine Translation difficult?

## Ambiguous Input

Kids make nutritious snacks

A million people live on water

## Variable Input

why is this complicated

why is this difficult

## Divergence between languages

word-order (syntactic structure), morphology, idiomatic expressions
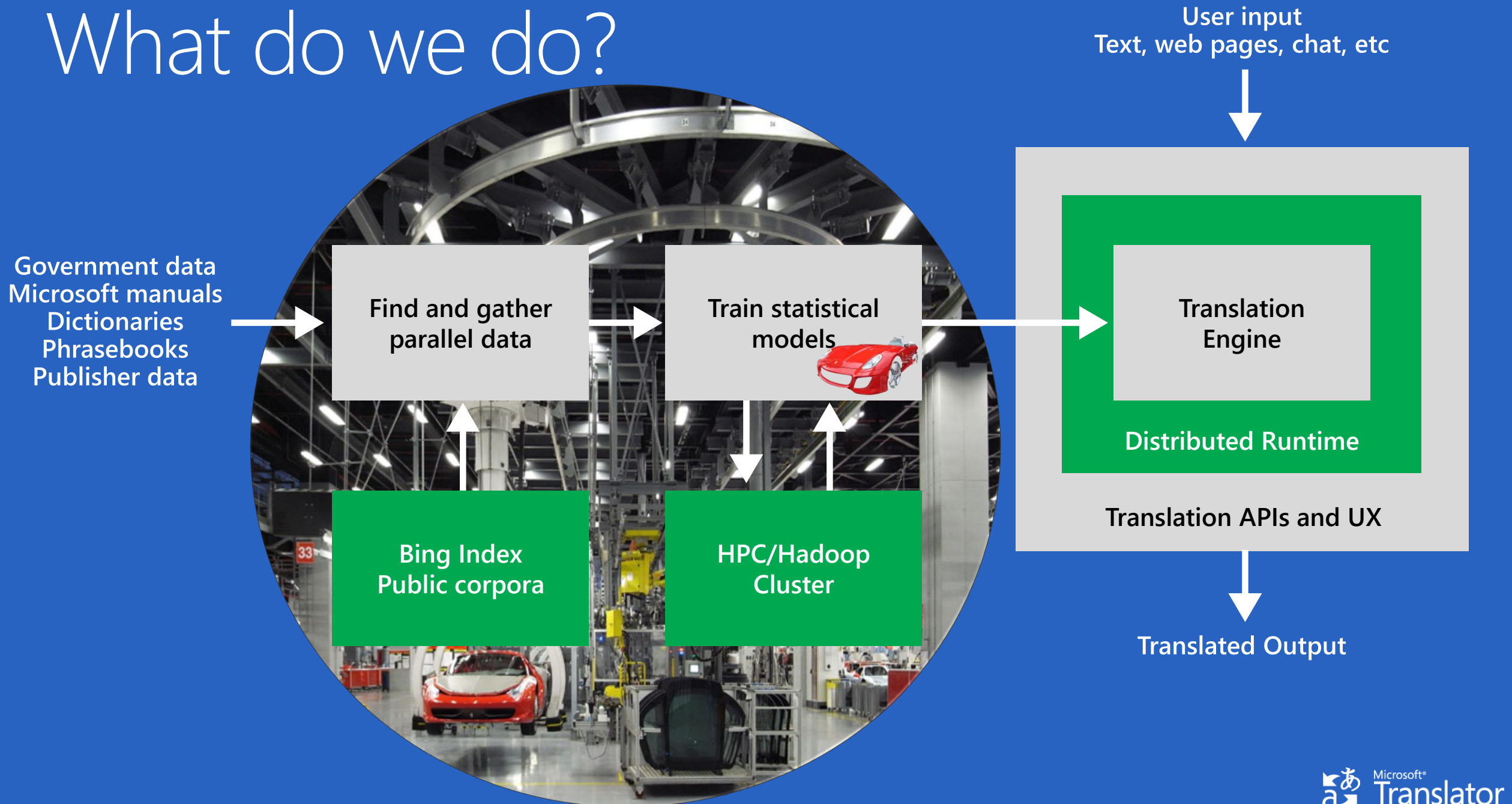
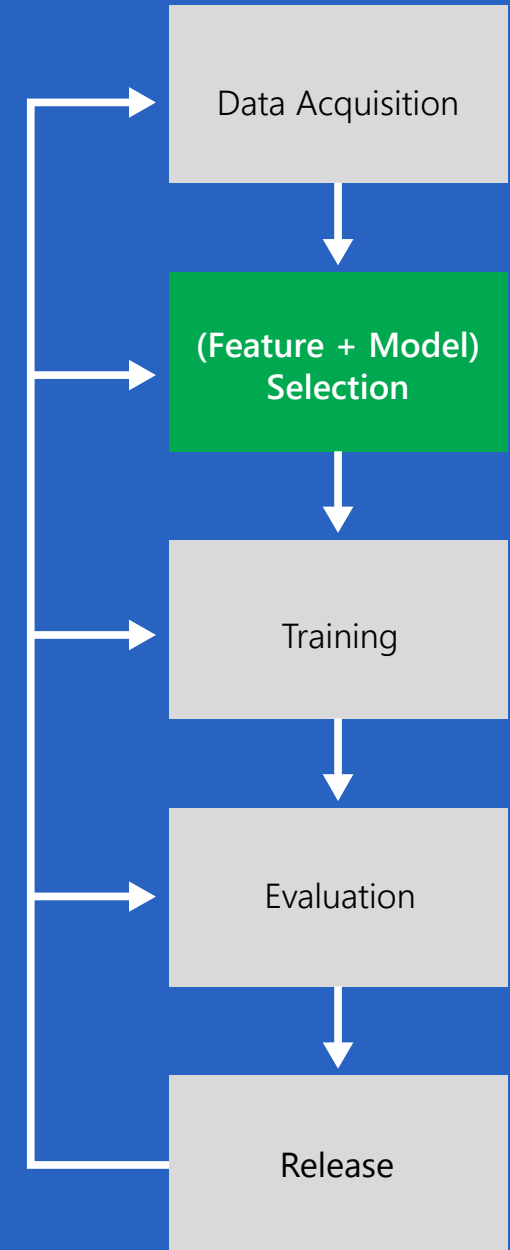## Evolving

just 4get abt it!

idk how did he reach b4 me

Microsoft® Translator
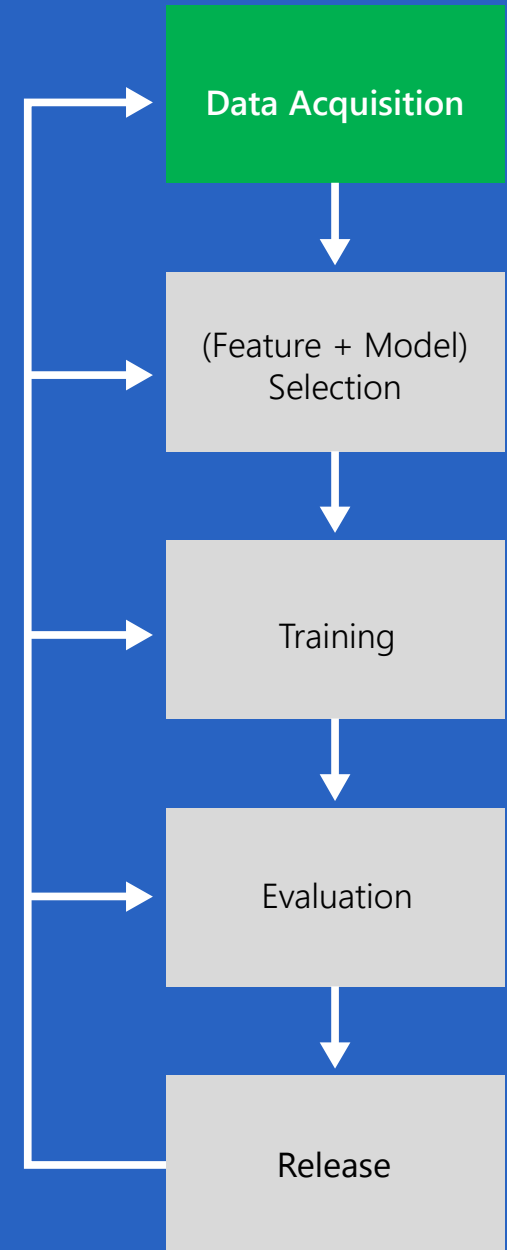
# Outline

Microsoft® Translator

# What do we do?

User input
Text, web pages, chat, etc

Government data
Microsoft manuals
Dictionaries
Phrasebooks
Publisher data

Find and gather
parallel data

Train statistical
models

Translation
Engine

Distributed Runtime

Bing Index
Public corpora

HPC/Hadoop
Cluster

Translation APIs and UX

Translated Output

Microsoft® Translator

# Outline: For-each

- Common Problems
- Our Solution
- Examples\Demo
- Summarize



Data Acquisition

(Feature + Model) Selection

Training

Evaluation

Release

Microsoft® Translator

# Data Acquisition

Data Acquisition

(Feature + Model) Selection

Training

Evaluation

Release

Microsoft® Translator

# Data Pipeline

Web

JRC. Europarl

Microsoft localization

→ Preprocess, Sentence Align → MT Store → Extract: Filter Chain → Train SMT Engine

Microsoft® Translator

# Pre & Post processing data

UnicodeNormalization_Filter

Length_Filter

EscapeSequenceFilter

SentenceBreak_Filter

LanguageDetect_Filter

HtmlParser_Filter

LatinScriptRatio_Filter

WordCount_Filter

Dedupe_Filter

VocabularySaturation_Filter

# Data Explorer

# Data: Solutions

# Summary

- Order from chaos – store data & metadata
- Automate thy steps
- Invest heavily in data pre-processing
  - Aggregate data to reduce training time without sacrificing quality
  - Junk data removal
- Once again, great data is the key! Both quantity and quality matter!

Microsoft® Translator

# Training: Common Issues



Easy, Quick?
Reproducible?
Stale?
Scaling?

# Push button training

# Monitoring

# Sample Config

```
<ProjectConfig>

    <SourceLanguage>ENU</SourceLanguage>

    <TargetLanguage>HIN</TargetLanguage>

    <ParallelSentencePath>\\mt-
data\data\datasink\mtmain_backend\paralleldata\general\en-us\hi-
in\train</ParallelSentencePath>

    <LambdaParallelSentencePath>\\mt-
data\data\datasink\mtmain_backend\paralleldata\general\en-us\hi-
in\dev</LambdaTrainSourceSentences>

    <SmokeTestParallelSentencePath>\\mt-
data\data\datasink\mtmain_backend\paralleldata\general\en-us\hi-
in\test</LambdaTrainSourceSentences>

</ProjectConfig>
```

# Expanded Config

```xml
<?xml version="1.0" encoding="utf-16"?>
<ProjectConfig xmlns:xsd="http://www.w3.org/2001/XMLSchema" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance">
  <DomainName>general</DomainName>
  <TrainMini>false</TrainMini>
  <SourceLanguage>ENU</SourceLanguage>
  <TargetLanguage>HIN</TargetLanguage>
  <EmailNotifyList>mttranno</EmailNotifyList>
  <ReferenceLengthSelection>Shortest</ReferenceLengthSelection>
  <SkipTranslationModelBuild>false</SkipTranslationModelBuild>
  <EuroparlBleu>false</EuroparlBleu>
  <DirectBleu>true</DirectBleu>
  <GenerateMobileModels>false</GenerateMobileModels>
  <MaxModelSize>50</MaxModelSize>
  <MobileModelsDropPath>\\mt-data\mobiletravel\mobiletravel\systems\latest_mobile</MobileModelsDropPath>
  <ParallelSentencePath>\\mt-data\data\datasink\mtmain_backend\paralleldata\general\en-us\hi-in\train</ParallelSentencePath>
  <ParallelSentencePathList>
    <ParallelSentencePath>\\mt-data\data\datasink\mtmain_backend\paralleldata\general\en-us\hi-in\train</ParallelSentencePath>
  </ParallelSentencePathList>
  <PreWordBrokenSrcDataPath />
  <PreWordBrokenTgtDataPath />
  <UseDataSelection>false</UseDataSelection>
  <GeneralData />
  <DomainData />
  <GenSelectData />
  <DomainDevData />
  <DataSelectionMethod>BiEntropyDiff</DataSelectionMethod>
  <DataSelectionRankMethod>DevDataThreshold</DataSelectionRankMethod>
  <DataSelectionLMOrder>4</DataSelectionLMOrder>
  <DataSelectionDiscount>0.7</DataSelectionDiscount>
  <SelectBestN>1000000</SelectBestN>
  <DataSelectionOutPath />
  <DataSelectionDryad />
  <DataSelectionThreshold>2</DataSelectionThreshold>
```

# Summary

- **Speed: Distributed, multi-threaded etc.**
  - As much effort into optimizing and scaling as runtime
  - Cache and re-use intermediate step results
- **Reliable training**
  - Every file, network access, etc
  - Auto resume – skip successful steps
- **Simplify E2E training**
  - Deterministic random seeds – data/code change?
  - Intelligent defaults for params, reduce human error
  - Continuously train, prevent stuff getting stale
  - Only E2E evals count for final shipping!

Microsoft® Translator

# Debugging & Evaluation

# Evaluation: Common Issues

# Evaluation





Microsoft® Translator

# Debugging: Common Issues

# Measure everything that matters

## Model Files

| Name | Size | Date Modified |
|---|---|---|
| enu.cym.general.back.dtable | 16,310 KB | 2/13/2014 2:56:38 PM |
| enu.cym.general.back.ttable | 232,332 KB | 2/13/2014 2:56:38 PM |
| enu.cym.general.fore.dtable | 14,342 KB | 2/13/2014 2:55:21 PM |
| enu.cym.general.fore.ttable | 232,272 KB | 2/13/2014 2:55:21 PM |
| enu.cym.general.lambdas.cmn | 2 KB | 2/13/2014 6:44:18 PM |
| enu.cym.general.mappingtable.dat64 | 201,720 KB | 2/13/2014 5:49:09 PM |
| enu.cym.general.mappingtable.idx64 | 65,030 KB | 2/13/2014 5:52:51 PM |
| enu.cym.general.mappingtable.mmlexmap64 | 7,932 KB | 2/13/2014 5:50:27 PM |
| enu.cym.general.noninteractive.cmn | 5 KB | 2/13/2014 6:45:41 PM |
| enu.cym.general.ordertemplate.mmtt64 | 10,000 KB | 2/13/2014 4:29:49 PM |

## Perf Results (Best)

| Corpus | Average | STDEV | Average99 | STDEV99 | Average95 | STDEV95 |
|---|---|---|---|---|---|---|
| gold_me | 0.0320 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| lambda_me | 0.1573 | 0.2181 | 0.1420 | 0.1567 | 0.1191 | 0.0686 |
| Profion-PAC-test | 0.1352 | 0.1058 | 0.1287 | 0.0931 | 0.1175 | 0.0801 |
| smoke_me | 0.1789 | 0.2809 | 0.1587 | 0.1165 | 0.1443 | 0.0923 |
| WMT_Test_000000000 | 0.1653 | 0.3330 | 0.1370 | 0.1331 | 0.1186 | 0.0685 |

## Perf Results (Worst)

| Corpus | Average1 | STDEV1 | Average5 | STDEV5 | Average10 | STDEV10 |
|---|---|---|---|---|---|---|
| gold_me | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| lambda_me | 1.7000 | 0.1266 | 0.9015 | 0.5569 | 0.5888 | 0.4962 |
| Profion-PAC-test | 0.0000 | 0.0000 | 0.4263 | 0.0848 | 0.3718 | 0.0817 |
| smoke_me | 2.1860 | 1.5942 | 0.8376 | 0.9816 | 0.5961 | 0.7342 |
| WMT_Test_000000000 | 2.9110 | 1.3447 | 1.0587 | 1.1458 | 0.6667 | 0.8950 |

Mon 1/13/2014 10:56 AM

**nlpsrv@microsoft.com**

MT ENU-ESN General - Nightly Smoke Test (mtmain_backend)

To  MT Nightly Development

If there are problems with how this message is displayed, click here to view it in a web browser.

Who's Who

### MT ENU-ESN General (Treelet) - Nightly Smoke mtmain_backend (Smoke Test/Lambda Only)

| Label | test match_alt_cat after Caboom fix |
|---|---|
| BLEU (smoke) | 35.31     (00.00) |

baseline  last

| OOV | 739 |
|---|---|
| Perf | Average for 95% = 0.1692 (00.54%) s |
| | 95 percentile value = 0.4680 s |
| | 99.5 percentile value = 0.8580 s |
| | 99.9 percentile value = 1.3110 s |

Performance Avg 95%

### History and Corpora

| Time | Build Number | smoke | lambda (train) | lambda (test) | fbmessages_xe | req_log_ex_no | s2s_swbd_test | s2s_appenstudio_ex | s2s_appenstudio_xe | wmt2009_test | fbmessages | mturkta_test776 | ted_lium_msra_utput |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 01/13/2014 18:39:53 | 10.01.9713.0001 | 35.(00.0 31 0) | 29.(00.0 04 1) | 33.3(00.00 0) | 35.2(00.00 8) | 31.8(- 8 00.01) | 41.48 (00.00) | 40.89 (00.00) | 28.4(00.0 60) | 32. 65(- 00.0 1) | 32.78(- 00.02) | 23.92 (00.00) |
| 01/08/2014 23:42:39 | 10.01.9707.0007 | 35. (00.0 31 0) | 29. (00.0 04 0) | 33.3(00.00 0) | 35.2(00.00 8) | 31.8(00.00 9) | 41.48 (00.00) | 40.89 (00.00) | 28.4(00.0 50) | 32. (00.0 66 0) | 32.80(00.00) | 23.92 (00.00) |
| 01/08/2014 03:19:39 | 10.01.9707.0007 | 35. (00.0 31 0) | 29. (00.0 04 0) | 33.3(00.00 0) | 35.2(00.00 8) | 31.8(00.00 9) | 41.48 (00.00) | 40.89 (00.00) | 28.4(00.0 50) | 32. (00.0 66 0) | 32.80(00.00) | 23.92 (00.00) |

Translator

# Debugging

```
The  author  wrote  the  novel  was  likely  to  be  a  best-seller.
DECL1————————NP1————————DETP1————————ADJ1*          "The"
                        NOUN1*        "author"
            VERB1*      "wrote"
            COMPCL1————NP2————————DETP2————————ADJ2*        "the"
                                  NOUN2*        "novel"
            VERB2*      "was"
            AJP1————————ADJ3*          "likely"
                        INFCL1————————INFTO1————————PREP1*        "to"
                        VERB3*          "be"
                        NP3————————DETP3————————ADJ4*          "a"
                                  NOUN3*          "best-seller"

CHAR1          "."
```

# Debugging

```
>display srgraph #332824
        analytics ( Cat"Noun" SubLinkId 1 Depth 1 AlignId 1 MatchSize 1)
        Links2> analytics ( Rule_no 10243 SubLinkId 1 Depth 1 AlignId 1 ChannelMLEFwdHeadLogp
            Prmods> google ( Depth 2 AlignId 1 ParentAttr"Prmods")
        analytique ( Rule_no 10244 SubLinkId 1 Depth 1 AlignId 1 ChannelMLEFwdHeadLo
        analyses ( Rule_no 10245 SubLinkId 1 Depth 1 AlignId 1 ChannelMLEFwdHeadLogp
        analytics (+DontLike Rule_no 10246 SubLinkId 1 Depth 1 AlignId 1 ChannelMLEF
        analyse ( Rule_no 10247 SubLinkId 1 Depth 1 AlignId 1 ChannelMLEFwdHeadLogpr
```

# Debugging Output



MTDiff - \\mt-data\NightlySmokeTests\NightlySmokeTests\mtmain_backend_core-xt\chs_enu_general\2013-11-06_06h21m02s\Current\bandb-eval.vs-previous.adf

File  Edit  View  Format  Bin  Tools  Help

| unsorted (487) | + the (3) | worse order (1) | better word choice (1) | equally good (1) | both bad (1) |

Edit Properties    Drag All Sentences    Delete

(5)
Src:      地处美丽的 Luberon 乡村边缘，距离 Avignon 半小时车程，40 分钟即可到达马赛机场。
Ref:      Half an hour from Avignon, forty minutes from Marseilles airport and on the edge of the beautiful Luberon countryside.
MT1/2:    generate2 {}
MT1:      Located in the beautiful Luberon village edge, drive from the Avignon for half an hour, 40 minutes to arrive at Marseille airport.
MT2:      Located in the edge of the beautiful Luberon village, drive from the Avignon for half an hour and 40 minutes to arrive at Marseille airport.
(8)
Src:      在这个带有大门的著名建筑内有电梯可以直达位于二层的公寓。
Ref:      There is lift access to this first floor apartment on this gated, prestigious development.
MT1/2:    generate2 {}
MT1:      In this famous with a gate building has elevator access to the second-story apartment.
MT2:      In the door of the famous building has elevator access to the second-story apartment.
(9)
Src:      公寓邻近法国一个顶级的 18 洞高尔夫球场，即将在未来几年开发新的 9 洞高尔夫球场，无论是业余爱好者还是专业高尔夫球手都会感到宾至如归。
Ref:      Situated around one of the leading 18 hole golf complexes in France, with a new 9 hole golf course to be developed in forthcoming years, both amateur and professi
MT1/2:    generate2 {}
MT1:      Close to France's top 18 holes golf courses, develops new 9-hole golf course in the coming years, both amateur and professional golfers will feel right at home.
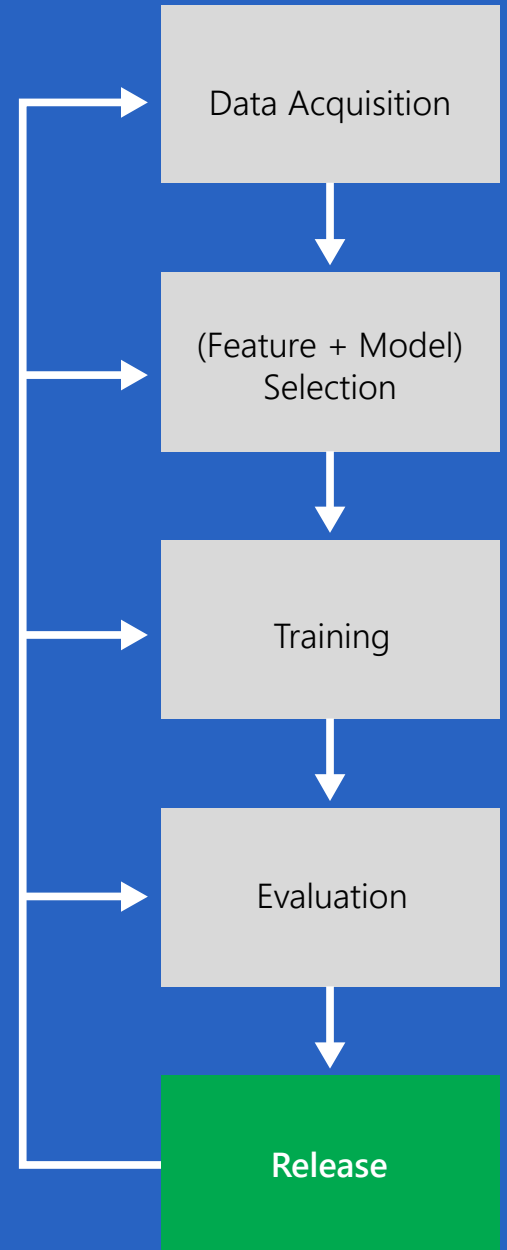MT2:      Close to France's top 18 holes golf course, will in the next few years to develop a new 9-hole golf course, both amateur and professional golfers will feel right at home

# Summary

- Log everything with E2E tracing

- Measure everything that matters
    - Aggregate vs. deterministic testing
    - Multiple test sets, refresh test sets
    - Performance, early indication
    - Tests for individual models as well as E2E composition

- Debugging
    - UI representation of decoding path
    - Configs for turning off non-determinism

Microsoft® Translator

# Common Pitfalls

# Summary

- Ship early and continuous measurement.
- Bridge expectation vs. Reality
- Models vs. code

vishalc@microsoft.com
@talktovishal

**Microsoft**

blogs.msdn.com/translator

twitter.com/MSTranslator

facebook.com/BingTranslator

Microsoft® Translator