

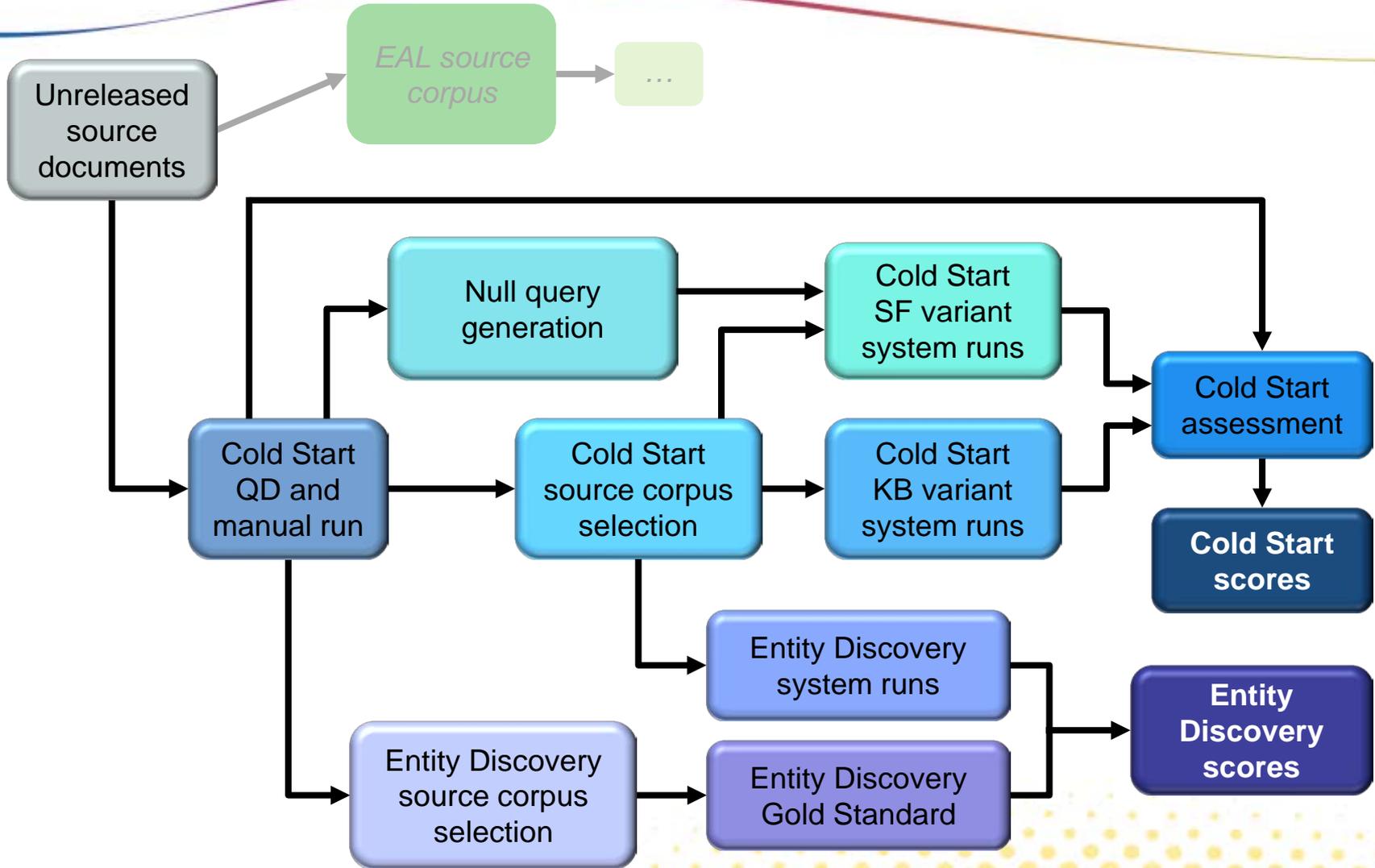


# **Linguistic Resources for the 2015 TAC KBP Cold Start and Tri-Lingual Entity Discovery & Linking Evaluations**

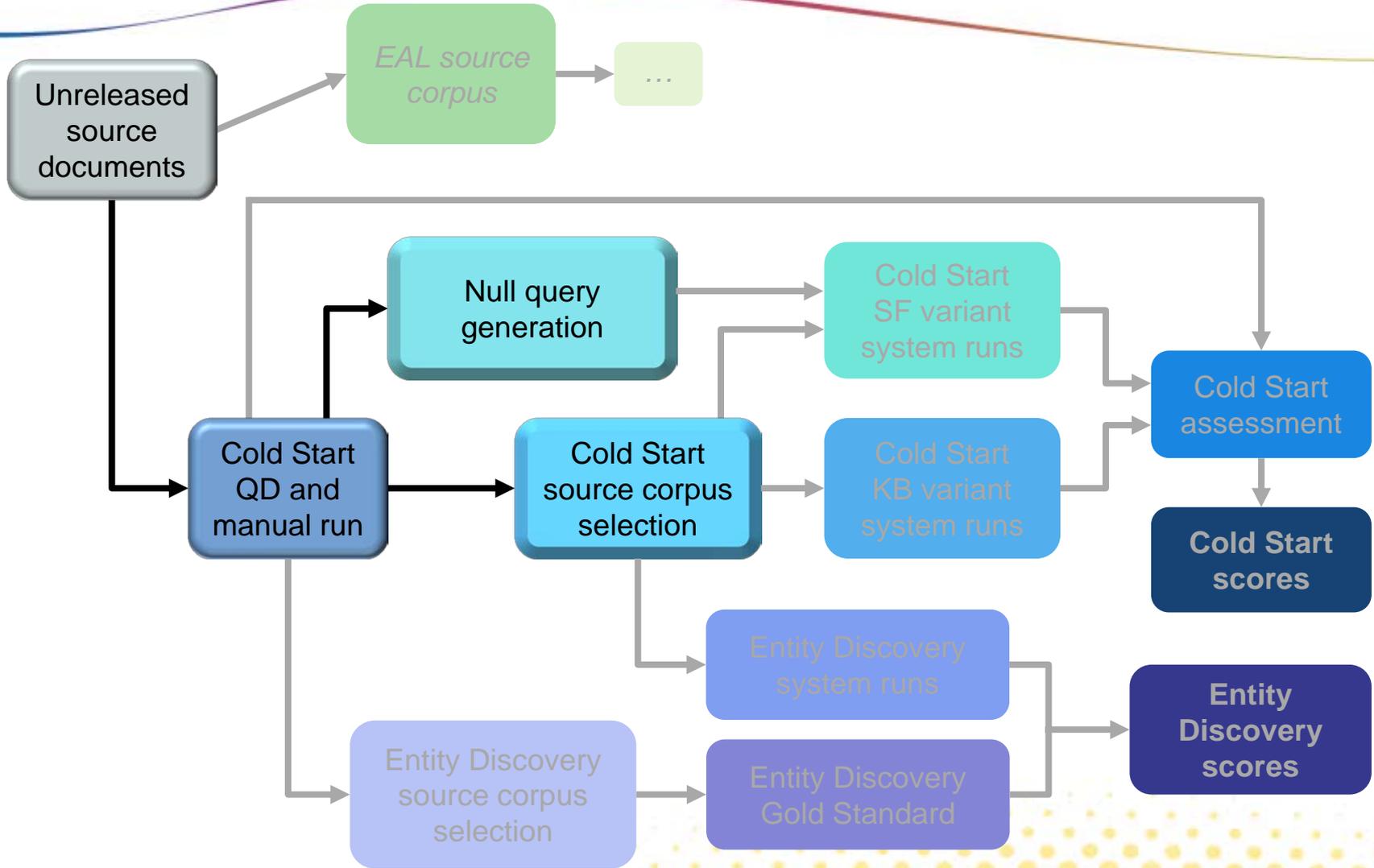
Joe Ellis (presenter), Jeremy Getman, Stephanie Strassel

*Linguistic Data Consortium  
University of Pennsylvania, USA*

- ◆ For data development purposes, Cold Start is a question answering task
  - Since 2012, LDC has approached Cold Start from the ‘Slot Filler Variation’ perspective
  - We’ve never previously had to concern ourselves much with the KB construction side of the task.
- ◆ However, query requirements changed significantly for 2015
  - More on this later...



- ◆ Three pools of unexposed documents
  - 2013 New York Times articles
    - ~57,000 documents
  - 2013 Xinhua articles
    - ~190,000 documents
  - Multi-post Discussion Forum threads (MPDFs)
    - Truncated discussion forum threads
    - Over 1 million MPDFs
- ◆ Annotators searched document pools to develop queries and the manual run
- ◆ Additional documents for the final source corpus also selected from these pools



Lance Barrett, 23, of London, KY, was charged with first-degree attempted burglary, theft of a firearm, and carrying a concealed weapon.

Lesa Bailey, 44, of London, KY, was charged with criminal conspiracy to make meth, unlawful possession of meth precursors and possession of a controlled substance.

London – *gpe:residents\_of\_city* – *per:charges*

• Lance Barrett

- first-degree attempted burglary
- theft of a firearm
- carrying a concealed weapon

• Lesa Bailey

- criminal conspiracy to make meth
- unlawful possession of meth precursors
- possession of a controlled substance

- ◆ Chains of entities connected by KBP slots
  - Cold Start queries comprised of
    - Entity – Slot 0 – Slot 1
- ◆ Cold Start annotation & query development concurrent
  - Annotators attempt to balance
    - Targeted number of annotations
    - Query variety (entity type, slot type, genre, etc.)
  - Annotation not exhaustive – some slots are more productive than others

# Cold Start: Query Development Changes

- ◆ Changes to query requirements compared to 2014 data
  - High degree of overlapping Entry Point Entities (EPEs) across queries
    - 2-5 mentions from different sources
    - Ambiguous whenever possible
  - Null queries
    - Auto-generated for rapid production but not guaranteed to have no valid responses
  - Changes made primarily to support Slot Filling subsumption and to ensure challenge for Entity Discovery

Search...

Type: None

Corpus: All Language: unreleased\_

Results

SourceDoc  
resulted, more women will be driven to clinics like this, which prosecutors called a "house of horrors."  
</P>  
<P>  
Last week, Judge Jeffrey P. Minehart of the Court of Common Pleas threw out three of seven first-degree murder charges against Gosnell. The doctor's defense lawyer, Jack J. McMahon, argued Monday that none of the remaining four cases had resulted in live births.  
</P>  
<P>  
Because the women were given injections of the drug digoxin, which causes "fetal demise," McMahon argued, any postdelivery movements were involuntary spasms.  
</P>  
<P>  
But Edward Cameron, an assistant district attorney, countered that testimony showed Gosnell did not always use digoxin and that it did not always work as intended. He quoted a former clinic worker with medical school training but no doctor's license who testified that the drug "wasn't giving the desired effect, the heart was always beating."  
</P>  
<P>  
Eight workers from the clinic, the Women's Medical Society in West Philadelphia, have pleaded guilty to lesser charges in the case, including Gosnell's wife, Pearl, a cosmetologist who helped perform abortions.  
</P>  
<P>  
If convicted, Gosnell could face the death penalty.  
</P>  
<P>  
Gosnell is also accused of third-degree murder in the death of a 41-year-old patient from Virginia, who visited his clinic after

QuerySelection Annotation

**Add Slot0** Expand/Collapse All

Slot0  
per:spouse

Filler0	Filler0 Entity Type	Normalization0
Pearl	PER	Specify...

Justification0  
Gosnell's wife, Pearl  
undefined  
undefined  
undefined

**Add Slot1**

Slot1  
per:employee\_or\_member\_of

Filler1	Filler1 Entity Type	Normalization1
Women's Medical Society	ORG	Specify...

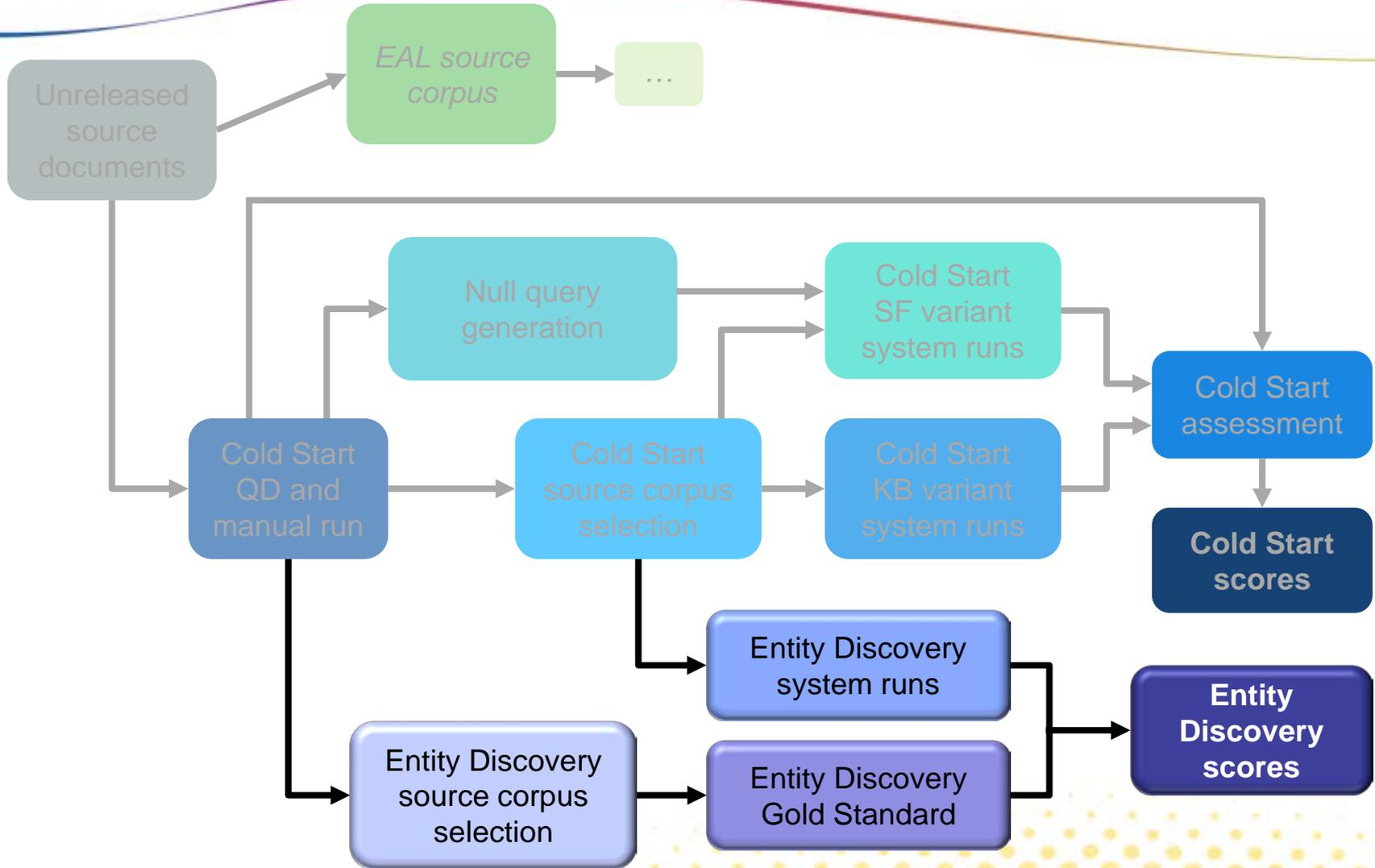
Justification1  
Eight workers from the clinic including Gosnell's wife,  
undefined  
undefined

**Delete Slot1**

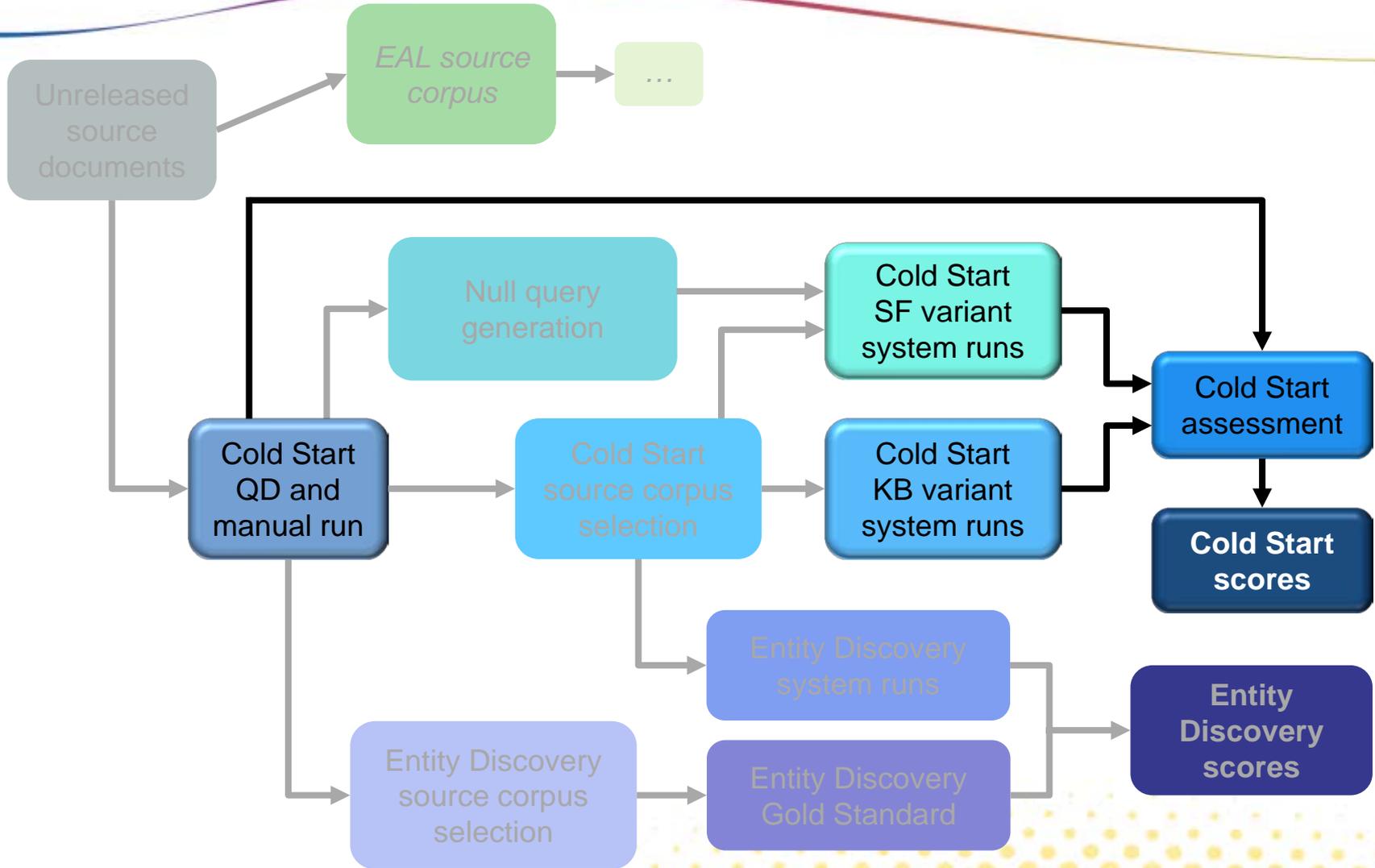
**Delete Slot0**

Slot0  
per:title

Filler0	Normalization0



- ◆ Identifying and clustering all valid entity types in the Cold Start corpus
  - Effectively, simplified Entity Discovery & Linking
    - One language, less entity types, one mention type
- ◆ Gold Standard development
  - As in ED&L, Cold Start – Entity Discovery submissions were scored against LDC’s gold standard mentions
  - 200 document subset of Cold Start source corpus
  - High degree of overlap with Cold Start queries and manual run
  - Entities mentioned in multiple documents, some with ambiguous mentions



- ◆ NIST pools results and sends to LDC
- ◆ Assessment performed in batches
  - hop-0 and hop-1 responses assessed for subset of queries
  - Queries to be assessed were selected and batched by NIST
- ◆ Assessment continues in batches until resources exhausted

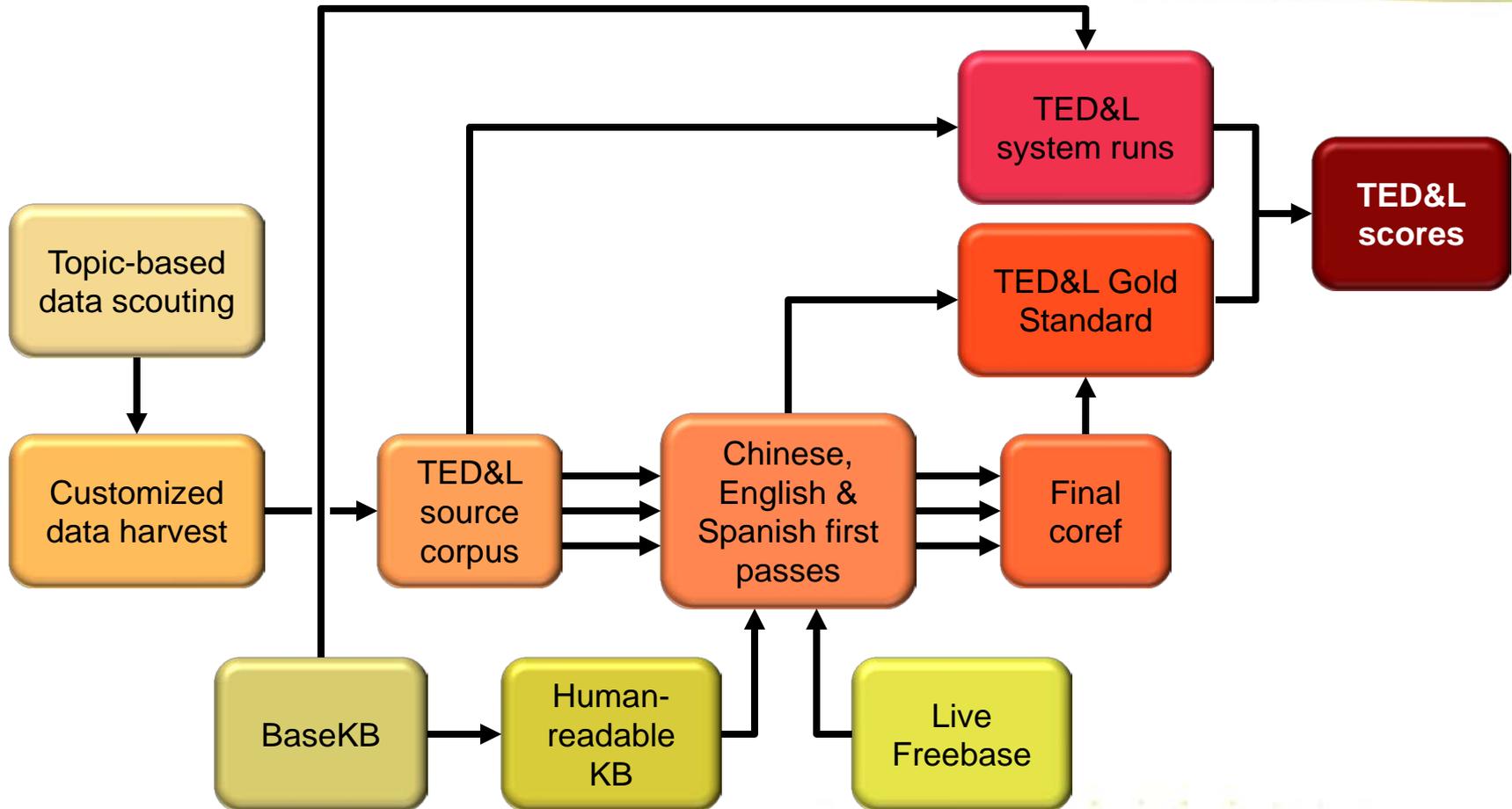
- ◆ Assess validity of fillers & justification from humans & systems
  - Filler
    - Correct – meets the slot requirements and supported in document
    - Wrong – doesn't meet slot requirements and/or not supported in doc
    - Inexact – otherwise correct, but is incomplete, includes extraneous text, or is not the most informative string in the document
  - Predicate
    - Correct – provides all information necessary to link the query entity to the filler by the chosen slot
    - Wrong – does not provide any of the necessary information
    - Inexact-Short – provides some, but not all, of the necessary information
    - Inexact-Long – otherwise correct, but includes extraneous text
- ◆ Correct and Inexact responses clustered together

# Cold Start: Assessment Results

	<b>Total</b>	<b>Newsire</b>	<b>MPDF</b>
Responses	30,654	15,948	14,706
Correct	26.7%	29.7%	23.5%
Wrong	68.8%	65.2%	72.8%
Inexact	4.5%	5.1%	3.7%

Track	Precision	Recall	F1
2014 Cold Start	91%	46%	62%
2015 Cold Start	81%	19%	31%

- ◆ New approach allowed for better tracking of query requirements, but may have further reduced focus on manual run
  - Focus given to competing query requirements
  - Annotators less exacting when selecting fillers
- ◆ Inexact responses included in scoring
- ◆ More queries
  - 1,327 productive queries (not including hop-1 portions)
  - 750 queries for Cold Start, Sentiment SF and Regular SF combined in 2014



- ◆ New knowledge base
- ◆ New source data requirements
- ◆ New annotation requirements
  - Monolingual to tri-lingual
  - New entity types
    - FAC & LOC
  - New mention type
    - NOM

- ◆ The old KB (2008 Wikipedia snapshot) made task too artificial
- ◆ Distributed via two releases
  - BaseKB (basekb.com)
    - FreeBase converted into RDF
  - Algorithm for creating for KB entries
    - Describes process by which triples were collected into pages for annotators to review

BaseKB	2008 Wikipedia Snapshot
As a triple store, systems can interact with the KB as a graph	Only available as a collection of entries
Over a billion facts about more than 40 million subjects	818K entries

## ◆ Requirements

- 500 documents
- Cross-lingual, cross-document entity mentions
- Multiple, varied, recent sources

## ◆ Challenges

- Unusual approach for harvesting/processing
  - Usual approach is larger volumes from fewer sources
  - Additional effort required
- Managing Intellectual Property Rights issues
  - Ensuring LDC has the right to annotate and redistribute collected data
  - 100s of sources required new approaches new approaches
    - Data distribution formats

<P>  
Dio, who was born Ronald James Padavona in Portsmouth, New Hampshire, and grew up in Cortland, New York, played in Black Sabbath in 1979. Ronnie's wife, Wendy, says he died on Sunday.  
<P>

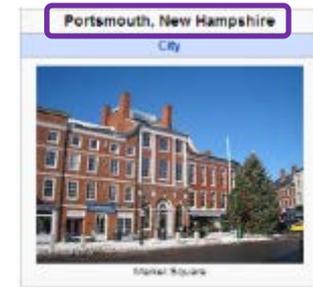


- ◆ Five entity types
  - PERs, ORGs, GPEs, FACs, LOCs
- ◆ Two mention types
  - Names and nominals
- ◆ Titles
  - Annotated to help distinguish PER nominal mentions
  - "the president [PER.NOM] signed a bill today"
  - "President[TTL] Clinton [PER.NAM] made a speech today"

## ◆ KB Linking

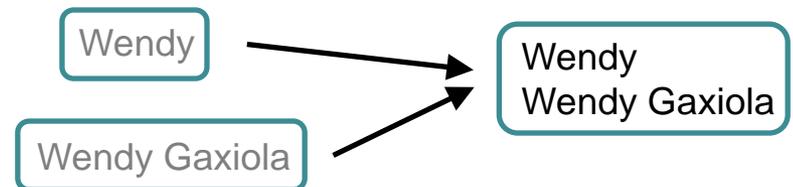
- Review ref document and search KB for matching node
- Multiple entities viewed together for quicker linking

<P>  
 Dio who was born Ronald James Padavona in Portsmouth New Hampshire, and grew up in Cortland, New York, played in Black Sabbath in 1979. His wife, Wendy, says he died on Sunday.  
 <P>



## ◆ NIL Coreference

- NIL queries (no KB match) require manual co-reference annotation



EntityDiscovery
EntityLinking

```
<?xml version="1.0" encoding="UTF-8"?>
<DOC id="ENG_NW_001001_20150719_F00100059" type="">
<HEADLINE>FIFA's Jeffrey Webb pleads not guilty in US
to corruption charges</HEADLINE>
<DATELINE/>
<TEXT>
<P>Jeffrey Webb, one of seven FIFA officials arrested
on corruption charges in Zurich, pleaded not guilty in his first
appearance in a U.S. court.</P>
<P>Webb, 50, who once served as FIFA's vice president,
was freed Saturday on $10 million bond and ordered by a U.S.
judge to remain under house arrest and around-the-clock guard.
</P>
<P>Unlike the six other officials at soccer's governing body
he was arrested with in May, Webb didn't fight extradition to
the U.S. to face trial, which may indicate that he's willing
to help prosecutors in an ongoing investigation of corruption
at FIFA. He was transferred to New York earlier in the week.
</P>
<P>His lawyer Edward O'Callaghan declined to comment after
court and wouldn't say if his client was cooperating with the
U.S.</P>
<P>If he does agree to help prosecutors, Webb could provide
```

Collapse/Expand All
Add Entity

Edward O'Callaghan

Entity Type:  
PER ▼

Entity Mentions:	Mention Type:	
Edward O'Callaghan	Name ▼	Delete
lawyer	Nominal ▼	Delete
lawyer	Nominal ▼	Delete
O'Callaghan	Name ▼	Delete
lawyer	Nominal ▼	Delete

KB Node ID:  
undefined

Web Search  Prose Search

NIL ▼

lawyer of Jeffrey Webb

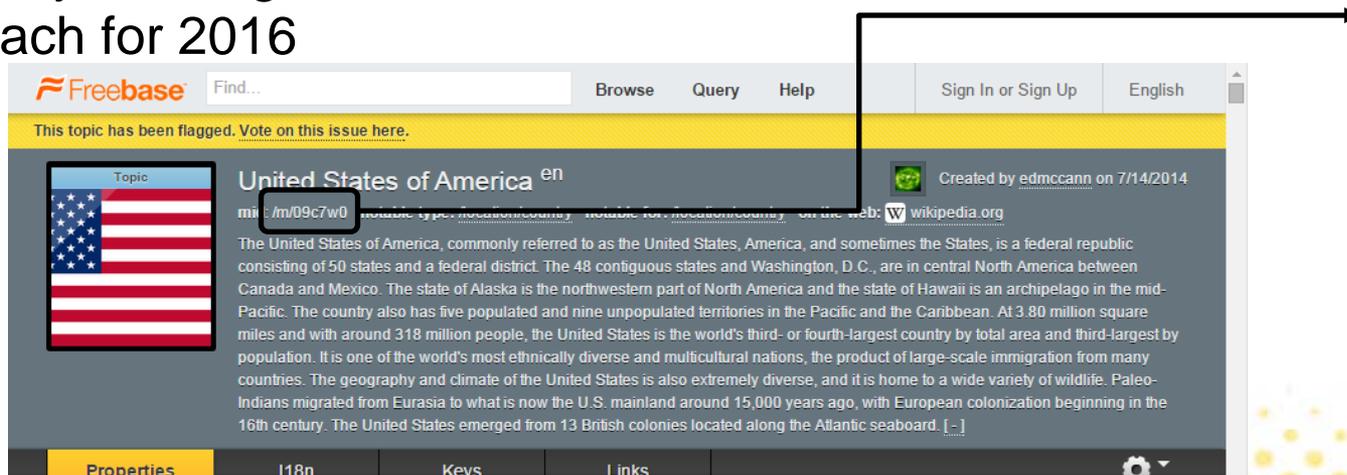
Delete Entity

Entity Type:  
TITLE ▼

Entity Mentions:	Mention Type:	
President		Delete

- ◆ Three types of annotators for within-doc kits
  - Monolingual English
  - Bilingual Chinese/English
  - Bilingual Spanish/English
- ◆ One English annotator for cross-doc coreference
  - Using English descriptions about the referents for each non-English cluster
    - Descriptions produced manually by original annotators during within-doc annotation

- ◆ Developed multiple indices for searching human-readable KB
  - None produced workable search results
    - Search for "united states" produced US page as 650<sup>th</sup> result
- ◆ Workaround
  - Annotators searched Freebase online and copied entity IDs
  - IDs not appearing in BaseKB converted to NILs during processing
- ◆ Currently working with NIST to create a better, more sustainable approach for 2016



	Training data (LDC2015E75)	Eval data (LDC2015E103)
<b>Total mentions</b>	<b>30,838</b>	<b>32,533</b>
ENG	13,545	15,645
CMN	13,116	11,066
SPA	4,177	5,822
<b>Total equivalence classes</b>	<b>5,744</b>	<b>7,235</b>
ENG	2,702	3,190
CMN	1,827	2,139
SPA	739	1,363
ENG/CMN	170	159
ENG/SPA	96	123
CMN/SPA	38	38
ENG/CMN/SPA	172	223

Catalog ID	Corpus Title	Size
<b>LDC2015E42</b>	TAC KBP 2015 Tri-lingual Entity Discovery and Linking Knowledge Base	1 knowledge base
<b>LDC2015E43</b>	TAC KBP 2015 Tri-lingual Entity Discovery and Linking Knowledge Base Entries Creation Algorithm	1 algorithm
<b>LDC2015E44</b>	TAC KBP 2015 Tri-lingual Entity Discovery and Linking Pilot Gold Standard Knowledge Base Links V1.1	686 mentions
<b>LDC2015E61</b>	TAC KBP 2015 Tri-lingual Entity Discovery and Linking Pilot Source Corpus	15 documents
<b>LDC2015E75</b>	TAC KBP 2015 Tri-lingual Entity Discovery and Linking Training Data V2.1	30838 mentions
<b>LDC2015E93</b>	TAC KBP 2015 Tri-lingual Entity Discovery and Linking Evaluation Source Corpus V2.0	500 documents
<b>LDC2015E102</b>	TAC KBP 2015 Tri-lingual Entity Discovery and Linking Evaluation Queries V1.2	32,533 queries
<b>LDC2015E103</b>	TAC KBP 2015 Tri-lingual Entity Discovery and Linking Evaluation Gold Standard Entity Mentions & Knowledge Base Links	32,533 mentions
<b>LDC2015E72</b>	TAC KBP 2015 English Cold Start Entity Discovery Sample Data	162 mentions
<b>LDC2015E76</b>	TAC KBP 2015 English Cold Start Evaluation Queries V2.0	2539 queries
<b>LDC2015E77</b>	TAC KBP 2015 English Cold Start Evaluation Source Corpus V2.0	49124 documents
<b>LDC2015E80</b>	TAC KBP 2015 English Cold Start Evaluation Queries and Manual Run	2218 responses
<b>LDC2015E81</b>	TAC KBP 2015 English Cold Start Entity Discovery Evaluation Gold Standard Entity Mentions V1.2	8110 mentions
<b>LDC2015E100</b>	TAC KBP 2015 English Cold Start Evaluation Assessment Results V3.1	30,678 assessments