

Introduction to Text Mining

Module 4: Development Lifecycle (Part 1)

Aims of this module

- Turning resources into applications: SLAM
- RichNews: multimedia application and demo
- Musing: Business Intelligence application
- KIM CORE Timelines application and demo
- GATE MIMIR: Semantic search and indexing in use
- The GATE Process

Semantic Annotation for the Life Sciences

Aim of the application

- Life science semantic annotation is much more than generic annotation of genes, proteins and diseases in text, in order to support search
- There are many highly use-case specific annotation requirements that demand a fresh look at how we drive annotation – our processes
- Processes to support annotation
 - Many use cases are ad-hoc and specialised
 - Clinical research – new requirements every day
 - How can we support this? What tools do we need?

Background

- The user
 - SLAM: South London and Maudsley NHS Trust
 - BRC: Biomedical Research Centre
 - CRIS: Case Register Information System
- February and March 2010
 - Proof of concept around MMSE
 - Requirements analysis, installation, adaptation
- Since 2010
 - In production
 - Cloud based system
 - Further use cases

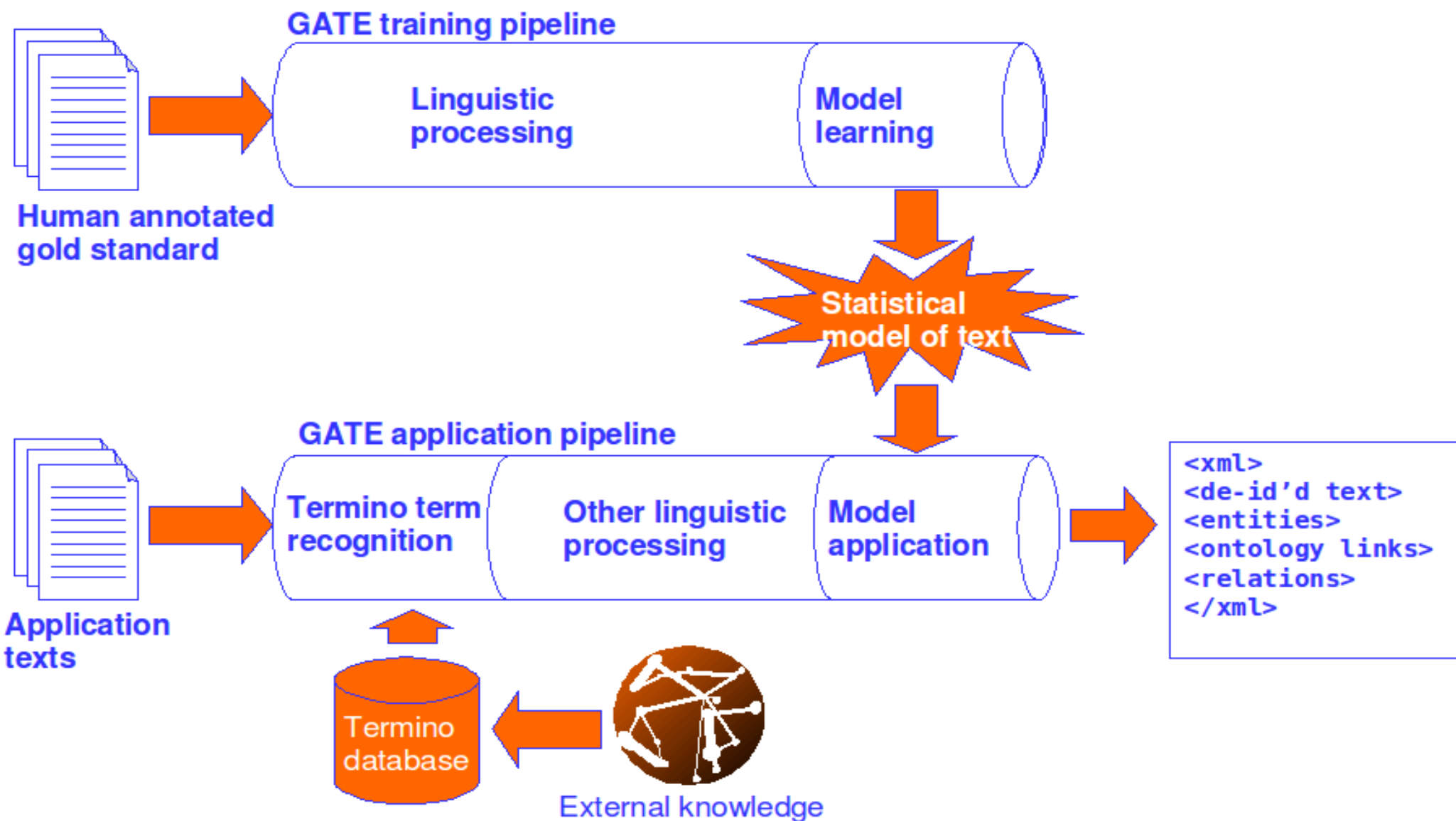
Clinical records

- Generic entities such as anatomical location, diagnosis and drug are sometimes of interest
- But many of the enquiries we have seen are more often interested in large numbers of very specific, and ad hoc entities or events
- This example is with a UK National Biomedical Research Centre
- An example – cognitive ability as shown by the MMSE score
- Illustrates a typical (but not the only) process

Types of IE systems

- Deep or shallow analysis
- Knowledge Engineering or Machine Learning approaches
 - Supervised
 - Unsupervised
 - Active learning
- GATE is agnostic

Supervised learning architecture





Messages GATE Corpus_001... 10004784-eEVCo1...

Annotation Sets Annotations List Annotations Stack Co-reference Editor Text



Unable to assess MMSE but last one on 1/1/8 was 21/30

assess MMSE but last one on 1/1/8 was 21/30. Bristol ADL score was 6/39 13/1/08; She currently

Type	Set	Start	End	Id	Features
MMSE	Automatic	983	1019	10846	{date=26/09/08, denominator=30, numerator=21, ruleMMSE=m
MMSE	forCorrection	983	1019	9860	{date=01/01/08, denominator=30, numerator=21, ruleMMSE=m

- MMSE-Lookup
- Number
- Score
- Sentence
- SpaceToken
- Split
- Token
- ▶ Original markups
- ▼ forCorrection
 - MMSE

2 Annotations (0 selected) Select:

New

Document Editor Initialisation Parameters



Messages GATE Corpus_001... 10096202-cATAt1...

Annotation Sets Annotations List Annotations Stack Co-reference Editor Text

CR8 2LJ Dear Dr Collins Re: ZZZZZ ZZZZZ DOB: ZZZZZ ZZZZZ ZZZZZ ZZZZZ ZZZZZ I reviewed Mrs. ZZZZZ at ZZZZZ on 6th March ZZZZZ in the presence of her daughter ZZZZZ and her son-in-law. Information was

Today she scored 5/30 on the MMSE. Care plan Mrs. ZZZZZ seems to be deriving little benefit from

- ▼ Automatic
 - Date
 - Lookup
 - MMSE
 - MMSE-Lookup
 - Number
 - Score
 - Sentence
 - SpaceToken
 - Split

Today she scored 5/30 on the MMSE

Type	Set	Start	End	Id	
MMSE	Automatic	926	942	1146	{date=09/03/09, denominator=30, numerator=5, ruleMMSE=score2}
MMSE	forCorrection	926	942	573	{date=06/03/09, denominator=30, numerator=5, ruleMMSE=score2}

2 Annotations (0 selected) Select:

Document Editor Initialisation Parameters

New



Messages GATE Corpus_001... 10096202-cATAt1...

Annotation Sets Annotations List Annotations Stack Co-reference Editor Text



CR8 2LJ Dear Dr Collins Re: ZZZZZ ZZZZZ DOR: 77777 77777

ZZZZZ on 6th March ZZZZZ in the presence of her daughter

Today she scored 5/30 on the MMSE. Care plan Mrs. ZZZZZ seems to be deriving little benefit from

I reviewed Mrs. ZZZZZ on 6th March

Number

Score

Sentence

SpaceToken

Split

Today she scored 5/30 on the MMSE

Type	Set	Start	End	Id	
MMSE	Automatic	926	942	1146	{date=09/03/09, denominator=30, numerator=5, ruleMMSE=score2}
MMSE	forCorrection	926	942	573	{date=06/03/09, denominator=30, numerator=5, ruleMMSE=score2}

2 Annotations (0 selected) Select:

New

Document Editor Initialisation Parameters

A shallow approach

- Pre-processing, including
 - morphological analysis
 - *“Patient was seen on” vs “I saw this patient on”*
 - POS tagging
 - *“patient was [VERB] on [DATE]”*
- Dictionary lookup
 - *“MMSE”, “Mini mental”, “Folstein”, “AMTS”*
- Coreference
 - *“We did an MMSE. It was 23/30”*

Annotations

His MMSE was 23/30 on 15 January 2008.

Annotations

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...


Annotations

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



Annotations

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



Id	Type
1	sentence

Annotations

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



Id	Type
1	sentence

Annotations

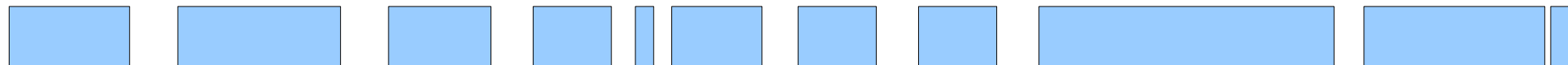
His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



Id	Type
1	sentence
2	token
3	token
4	token
5	token
6	token
7	token

Annotations

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



Id	Type	Start	End
1	sentence	0	39
2	token	0	3
3	token	4	8
4	token	9	12
5	token	13	15
6	token	15	16
7	token	16	18

Annotations

His MMSE was 23/30 on 15 January 2008.
 0...5...10...15...|...|...|...|...



Id	Type	Start	End	Features
1	sentence	0	39	
2	token	0	3	pos=PP
3	token	4	8	pos=NN
4	token	9	12	pos=VB
5	token	13	15	pos=CD
6	token	15	16	pos=SM
7	token	16	18	pos=CD

Annotations

His MMSE was 23/30 on 15 January 2008.
 0...5...10...15...|...|...|...|...



Id	Type	Start	End	Features
1	sentence	0	39	
2	token	0	3	pos=PP
3	token	4	8	pos=NN
4	token	9	12	pos=VB root=be
5	token	13	15	pos=CD type=num
6	token	15	16	pos=SM type=slash
7	token	16	18	pos=CD type=num

Dictionary lookup

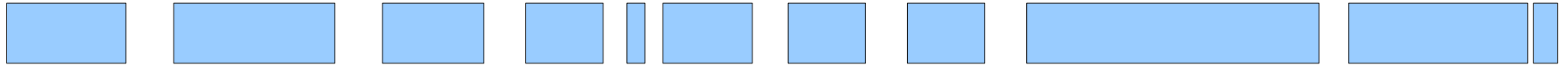
His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



Month

Dictionary lookup

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



MMSE

Month

Limitations of dictionary lookup

- Dictionary lookup is designed for finding simple, regular terms and features
- False positives
 - *“He may get better”*
 - *“Mother is a smoker”*
 - *“He often burns the toast, setting off the smoke alarm”*
- Cannot deal with complex patterns
 - For example, recognising e-mail addresses using just a dictionary would be impossible
- Cannot deal with ambiguity
 - I for Iodine, or I for me?

Pattern matching

- The early components in a GATE pipeline produce simple annotations (Token, Sentence, Dictionary lookups)
- These annotations have features (Token kind, part of speech, major type...)
- Patterns in these annotations and features can suggest more complex information
- We use JAPE, the pattern matching language in GATE, to find these patterns

Patterns

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



MMSE

Month

{number} {Month} {number}

Patterns

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



MMSE

Month

Date

Patterns

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



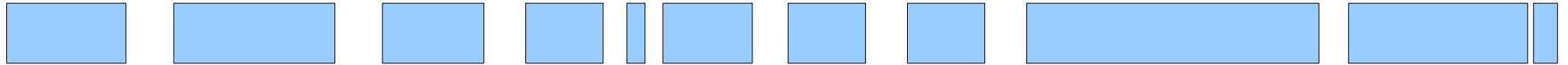
MMSE

Month

{number} {slash} {number}

Patterns

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



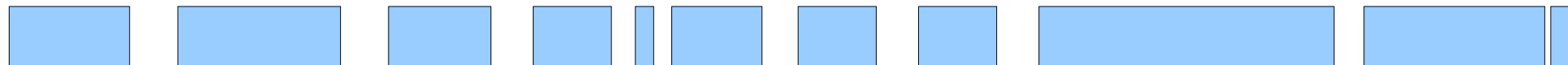
MMSE

Month

Score

Patterns

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



MMSE

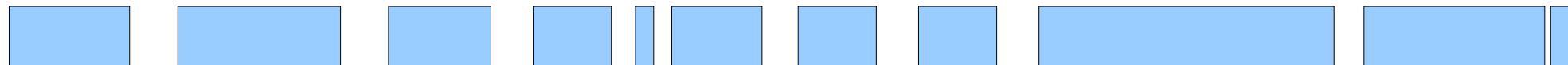
Month

Score

Date

Patterns

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



MMSE

Month

Score

Date

{MMSE} {BE} {Score} {?} {Date}

Patterns

His MMSE was 23/30 on 15 January 2008.
0...5...10...15...|...|...|...|...



MMSE

Month

Score

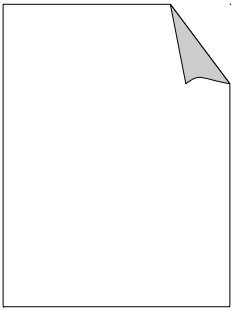
Date

MMSE with score and date

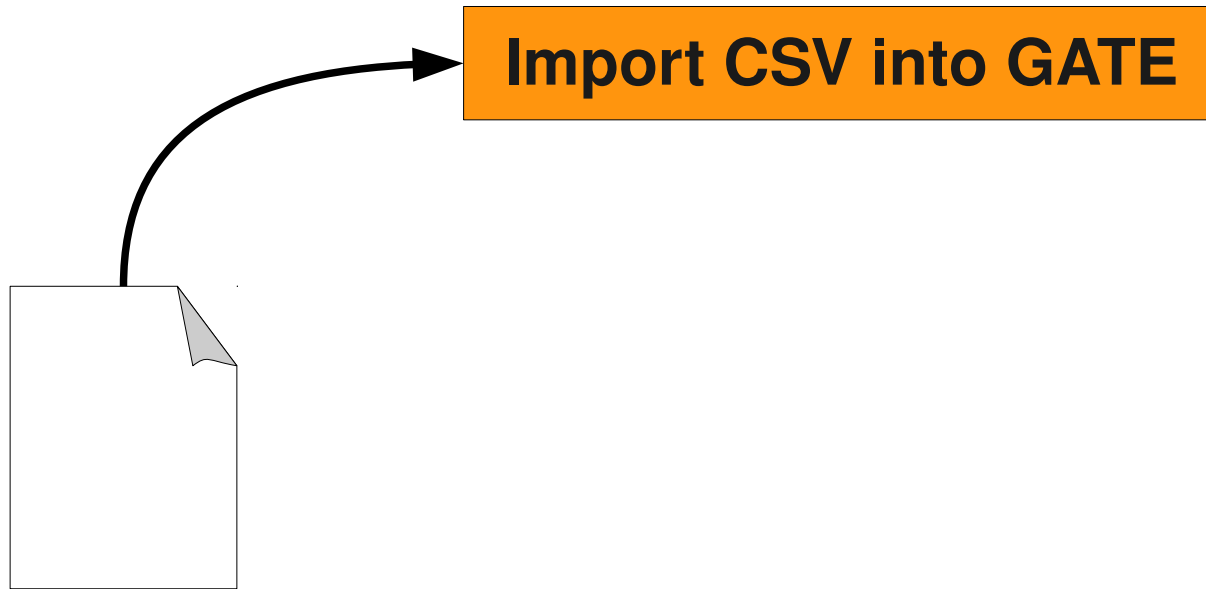
Patterns are general

- MMSE was 23/30 on 15 January 2009
- Mini mental was 25/30 on 12/08/07
- MMS was 25/30 last week
- MMSE is 25/30 today
- With adaptation
 - MMSE 25 out of 30
 - Long range dependencies on dates

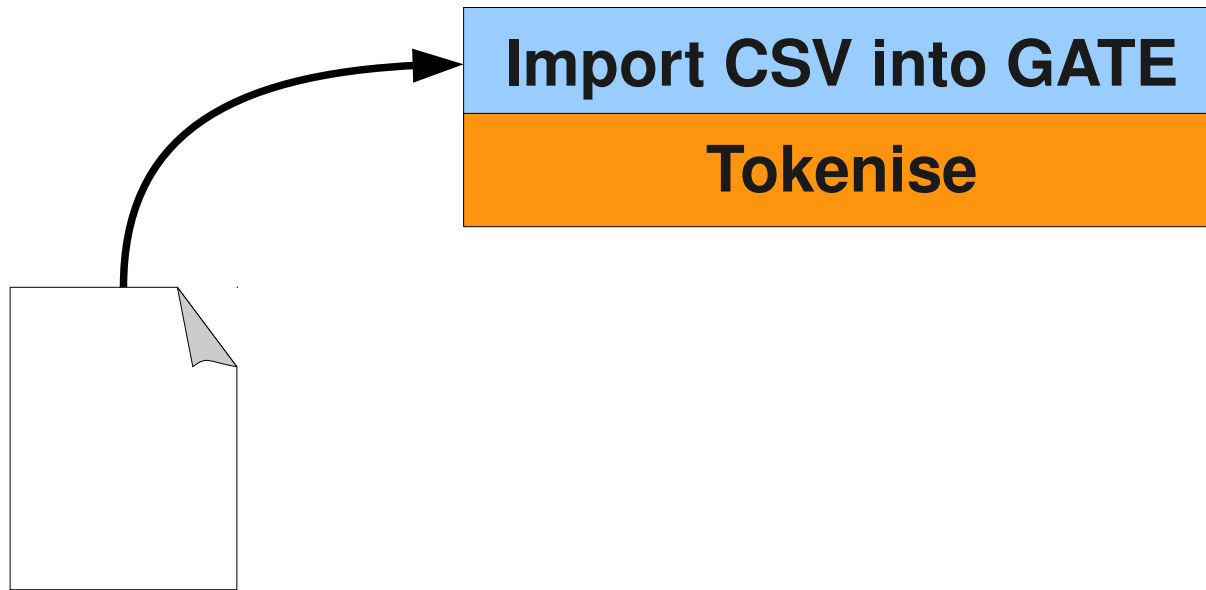
MMSE pipeline



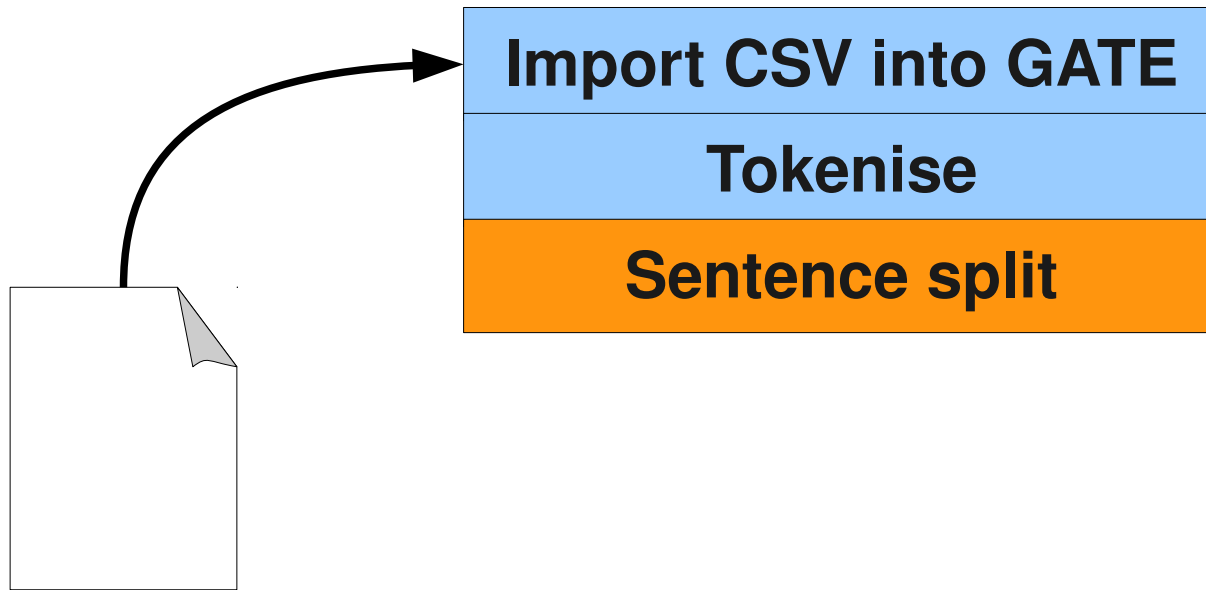
MMSE pipeline



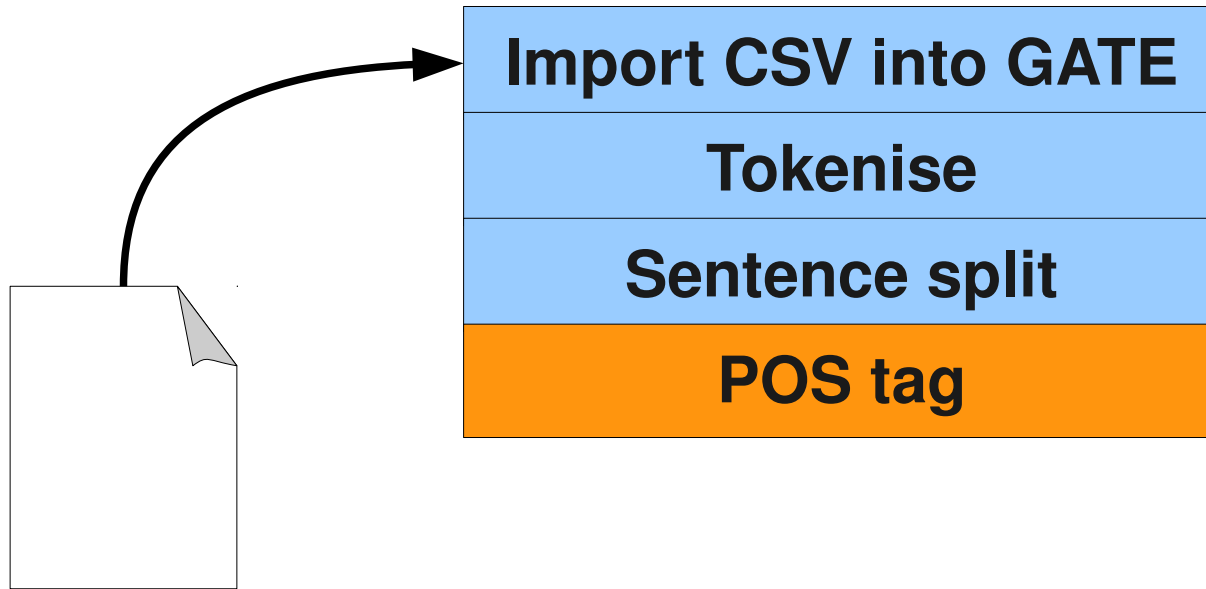
MMSE pipeline



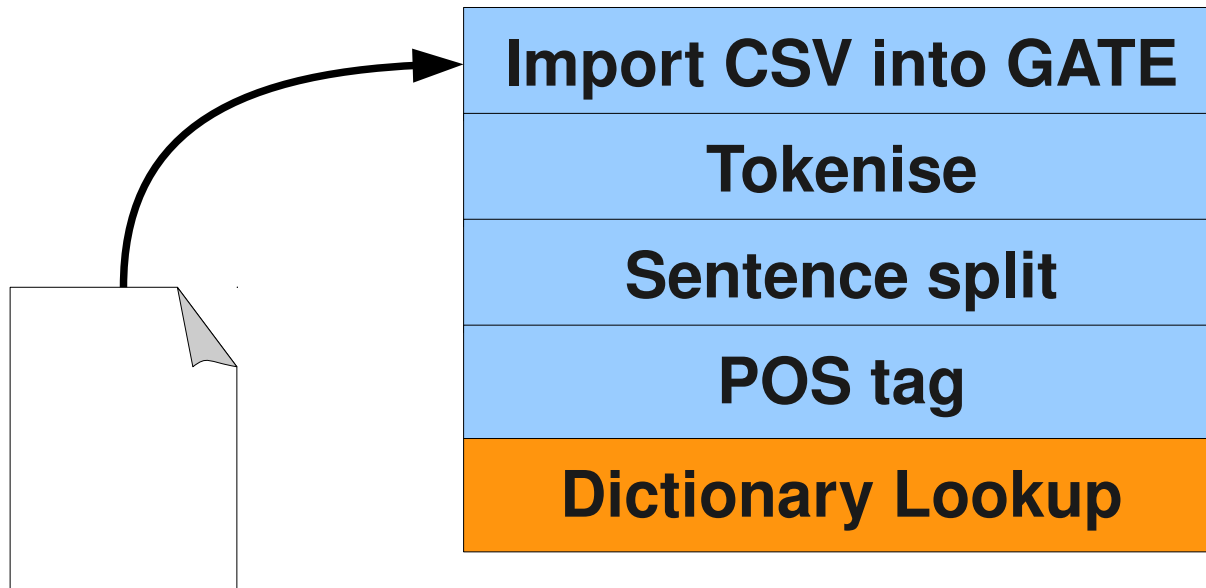
MMSE pipeline



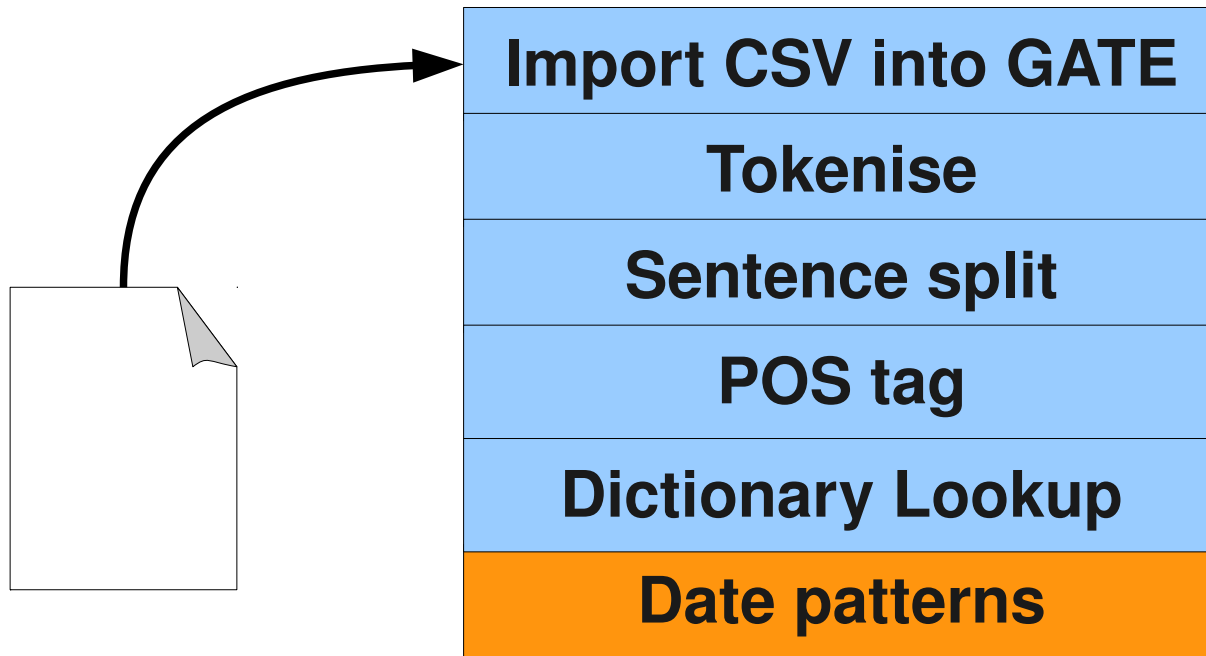
MMSE pipeline



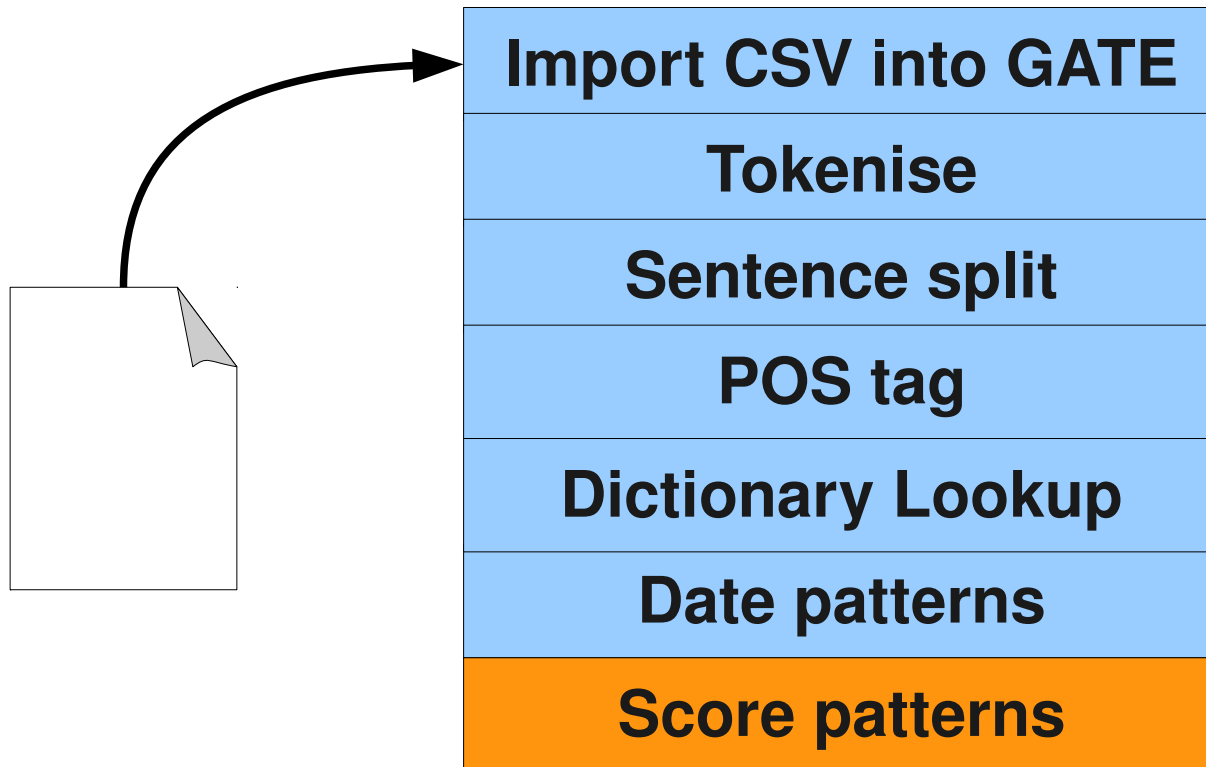
MMSE pipeline



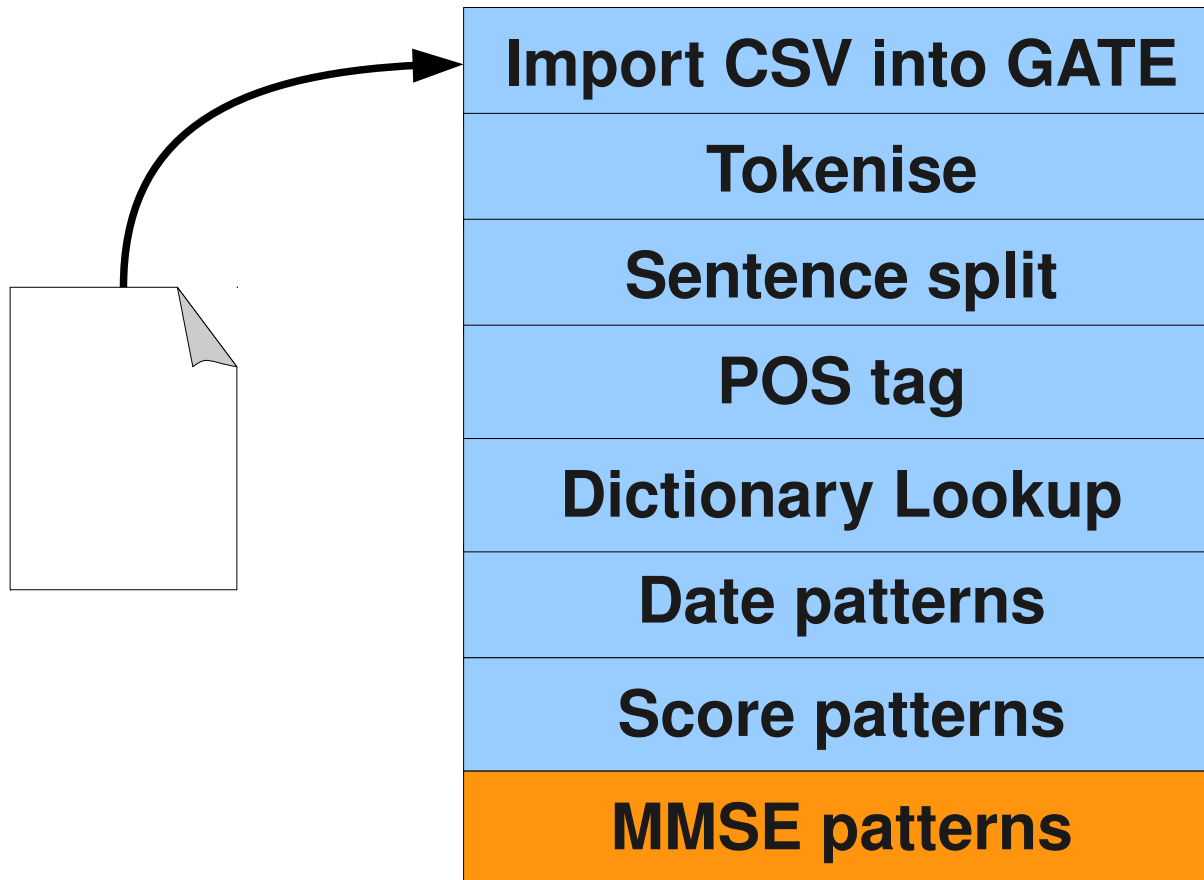
MMSE pipeline



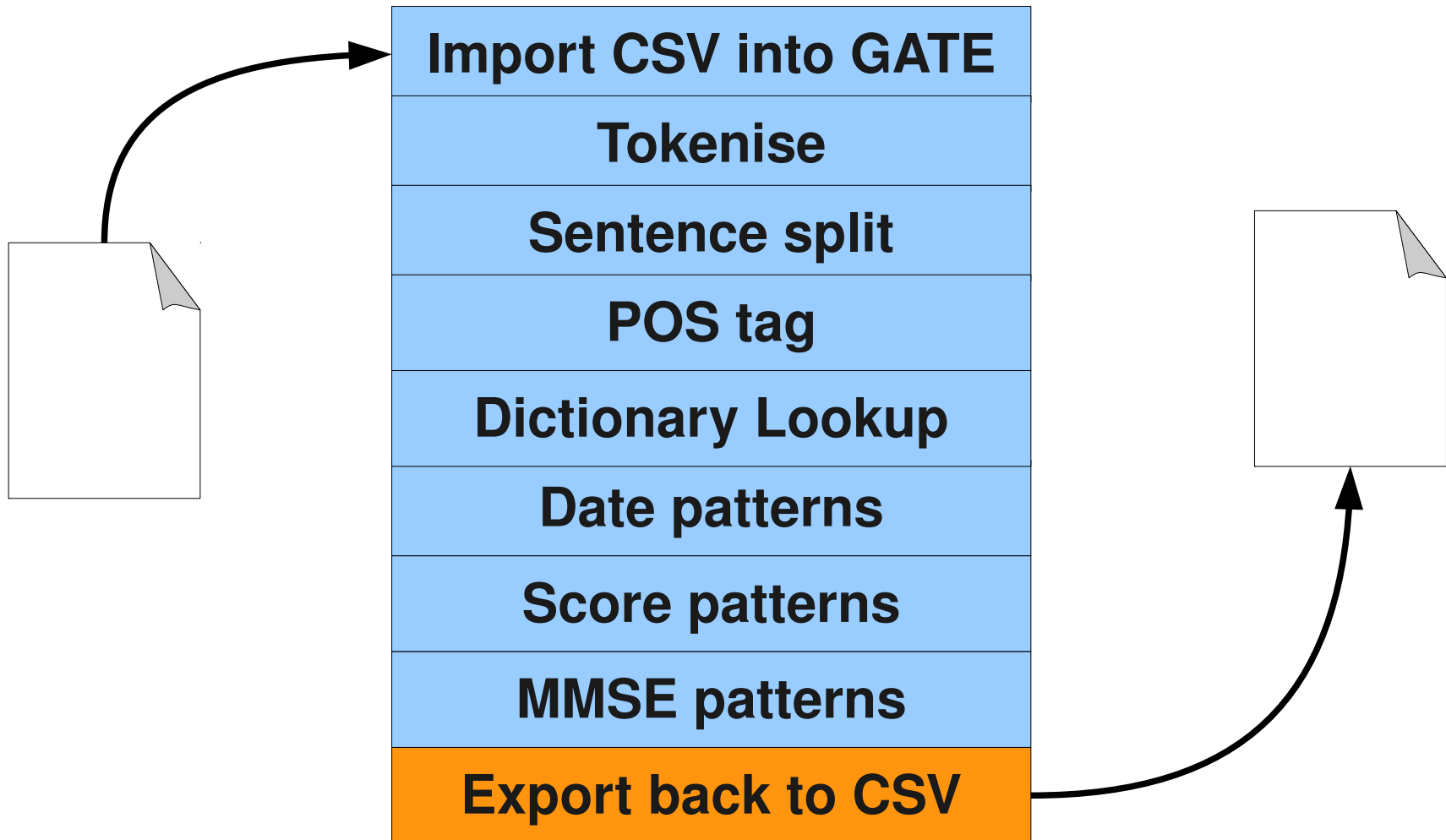
MMSE pipeline



MMSE pipeline



MMSE pipeline



Writing patterns

- Requires training
- Depending on time and skills, domain expert may take on some rule writing
- Requirements not always clear, and users do not always understand what the technology can do
- Needs a process to support
 - Domain expert manually annotates examples
 - Language engineer writes rules
 - Measure accuracy of rules
 - Repeat

The process as agile development

- IE system development is often linear
 - Guidelines → annotate → implement
- This is similar to the “waterfall” method of software development
 - Gather requirements → design → implement
- This has long been known to be problematic
- In contrast, our approach is **agile**

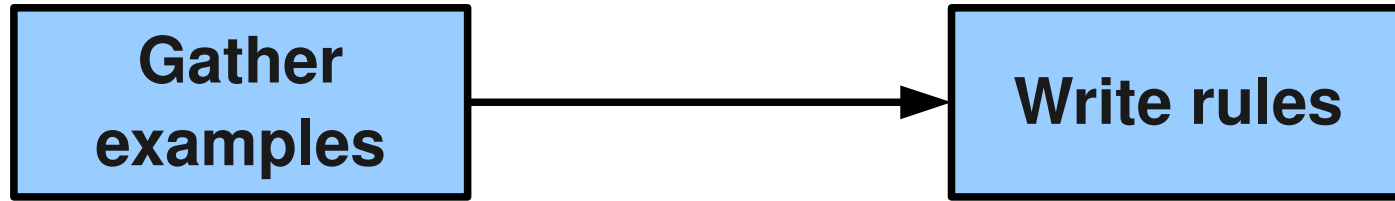
The process as agile development

- Recognise that requirements change
- Embrace that change
- Use it to drive development
- Developers and software engineers work alongside each other to understand requirements
- Early and iterative delivery
- Feedback to collect further requirements
- Reduces cost of annotation

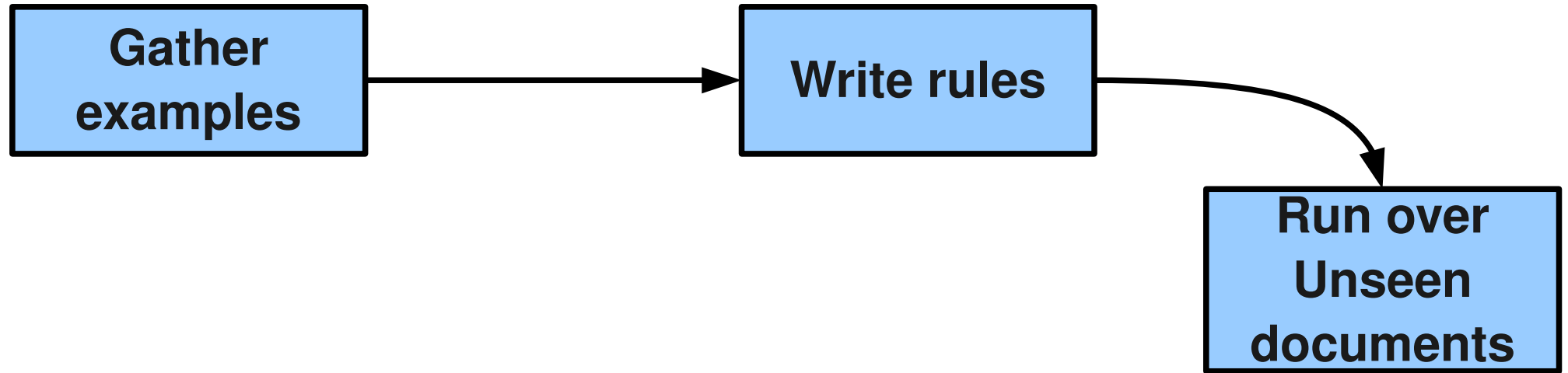
Annotation - a process

**Gather
examples**

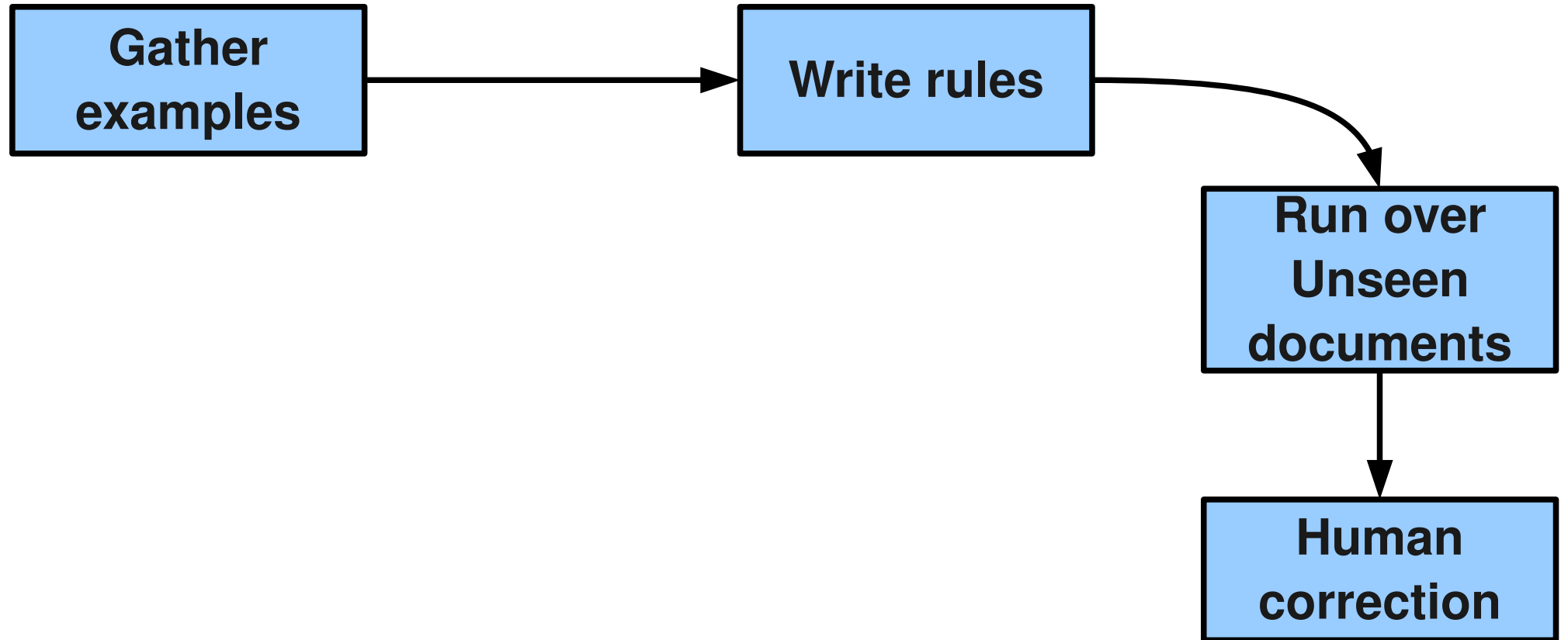
A process



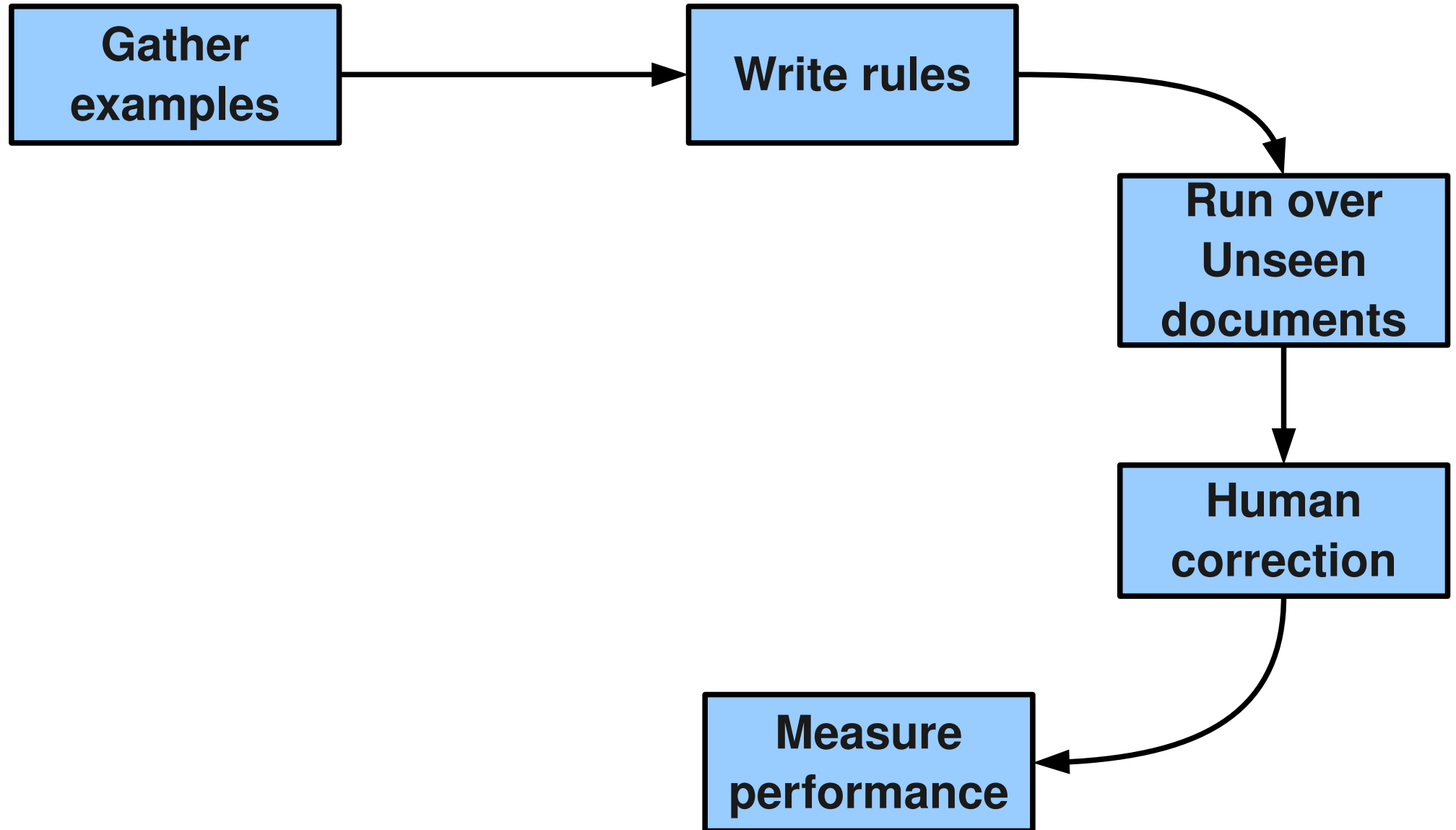
A process



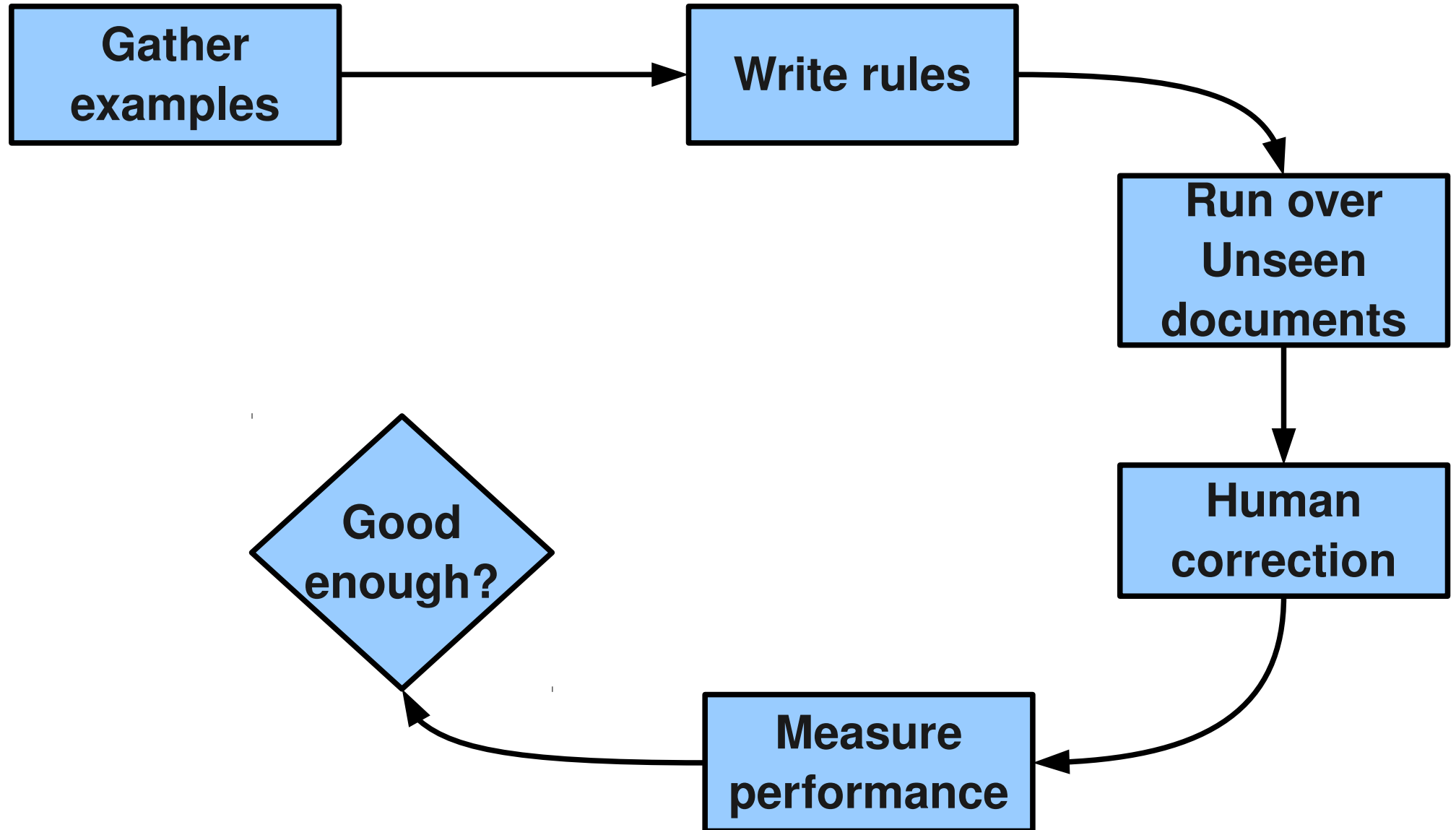
A process



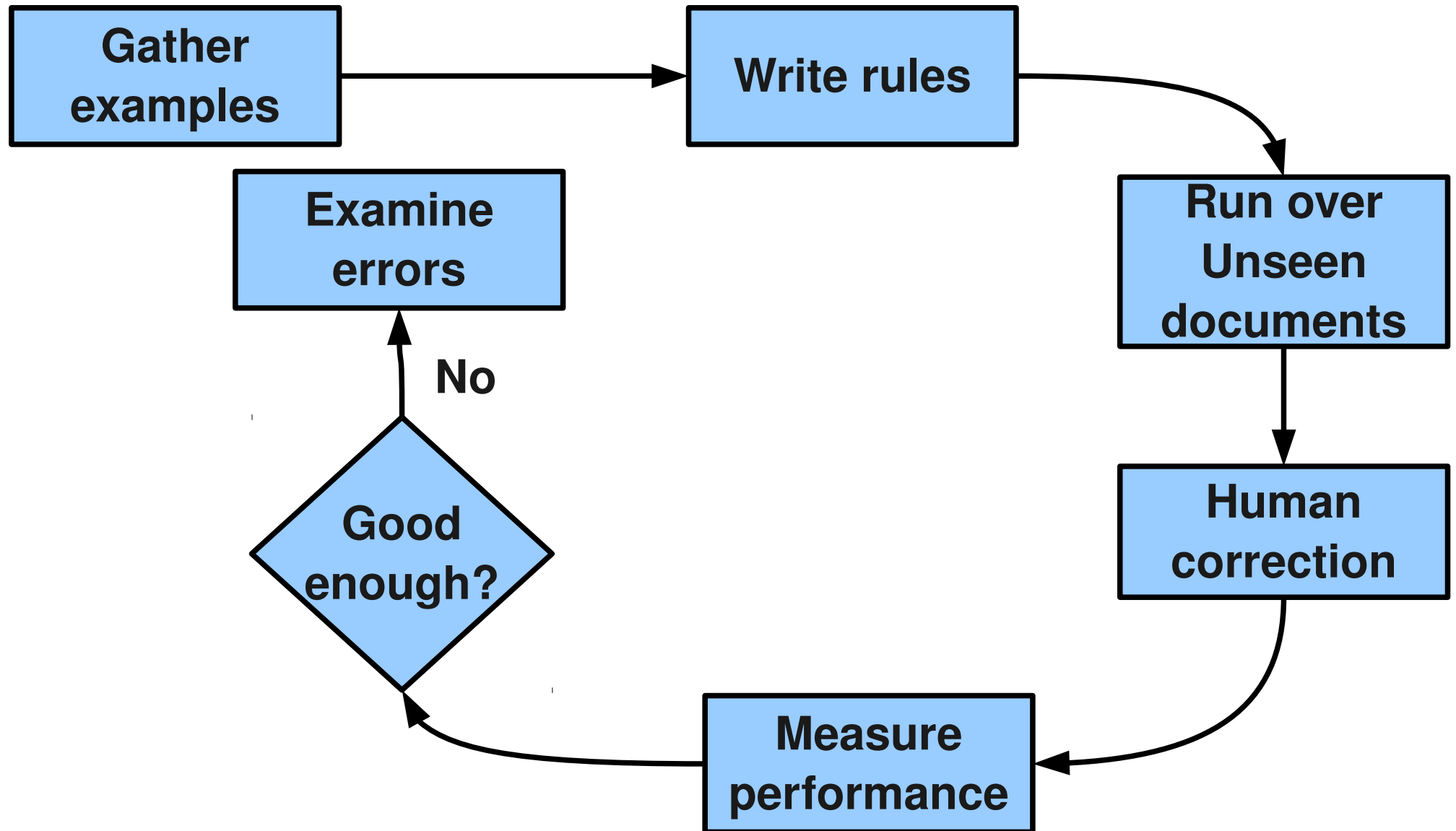
A process



A process



A process



Document	Match	Only A	Only B	Overlap	Rec. B/A	Prec. B/A	F1-lenient
-eEVCo18_3.xml_00103	1	0	0	0	1.00	1.00	1.00
-eEVCo18_4.xml_00104	1	0	0	0	1.00	1.00	1.00
-eEVCo18_0.xml_00105	2	0	0	0	1.00	1.00	1.00
-eEVCo18_1.xml_00106	1	0	0	0	1.00	1.00	1.00
-eEVCo18_2.xml_00107	1	0	0	0	1.00	1.00	1.00
-eEVCo18_0.xml_00108	1	0	0	0	1.00	1.00	1.00
-eEVCo18_3.xml_00109	1	0	0	0	1.00	1.00	1.00
-eEVCo18_3.xml_0010A	2	0	0	0	1.00	1.00	1.00
-eEVCo18_4.xml_0010B	0	1	1	0	0.00	0.00	0.00
-eEVCo18_0.xml_0010C	1	0	0	0	1.00	1.00	1.00
-eEVCo18_1.xml_0010D	1	0	0	0	1.00	1.00	1.00
-eEVCo18_3.xml_0010E	2	0	0	0	1.00	1.00	1.00
-eEVCo18_0.xml_0010F	1	0	0	0	1.00	1.00	1.00
-eEVCo18_2.xml_00110	0	0	2	0	1.00	0.00	0.00
-eEVCo18_12.xml_00111	1	0	0	0	1.00	1.00	1.00
-eEVCo18_13.xml_00112	1	0	0	0	1.00	1.00	1.00
-eEVCo18_14.xml_00113	1	0	0	0	1.00	1.00	1.00
-eEVCo18_16.xml_00114	1	1	1	0	0.50	0.50	0.50
-eEVCo18_17.xml_00115	1	0	0	0	1.00	1.00	1.00

Annotation Diff

Annotation Diff Tool

Document: ft-bank-of-england.xml_00016
 Annotation Set: Key
 Annotation Type: Date
 F-Measure Weight: 1.00

Response: ft-bank-of-england.xml_00016
 Annotation Set: [Default set]
 Features: All Some None

Do Diff

Start	End	Key	Features	=?	Start	End	Response	Features
1318	1332	second quarter	{kind= date}	-?				
1466	1474	Thursday	{}	-?				
212	222	early 1964	{kind= date}	~	218	222	1964	{kind= date, rule1=TempYear3, rule2=YearOnlyFinal}
23	31	Thursday	{kind= date, rule1=GazDate, rule2=DateOnlyFinal}	=	23	31	Thursday	{kind= date, rule1=GazDate, rule2=DateOnlyFinal}
1005	1015	last month	{kind= date}	=	1005	1015	last month	{kind= date, rule1=ModifierDate, rule2=DateOnlyFinal}
1582	1591	next week	{kind= date}	=	1582	1591	next week	{kind= date, rule1=ModifierDate, rule2=DateOnlyFinal}

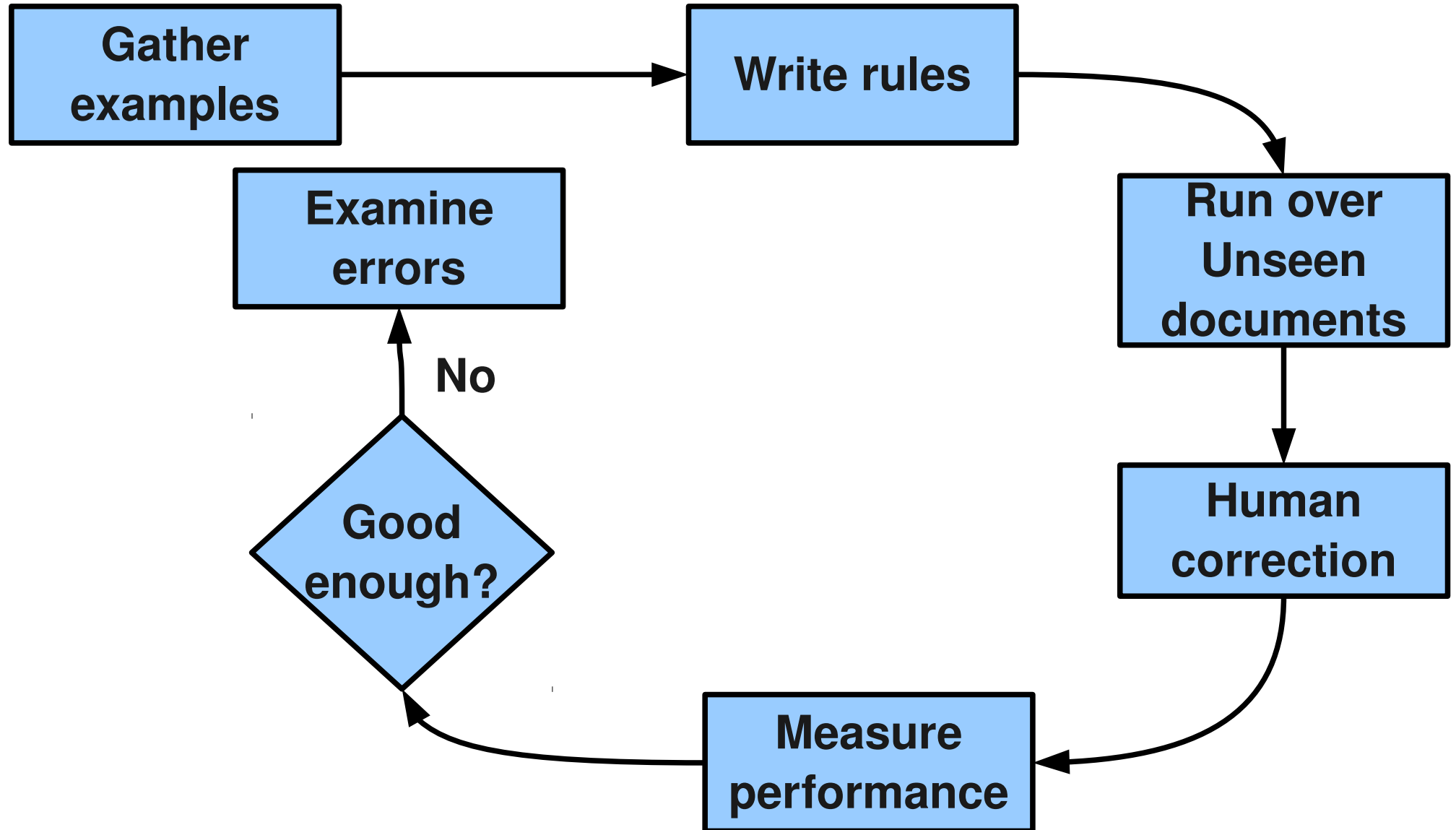
Correct: 3 Recall Precision F-Measure **Export to HTML**

Partially Correct: 1 Strict: 0.50 0.75 0.60

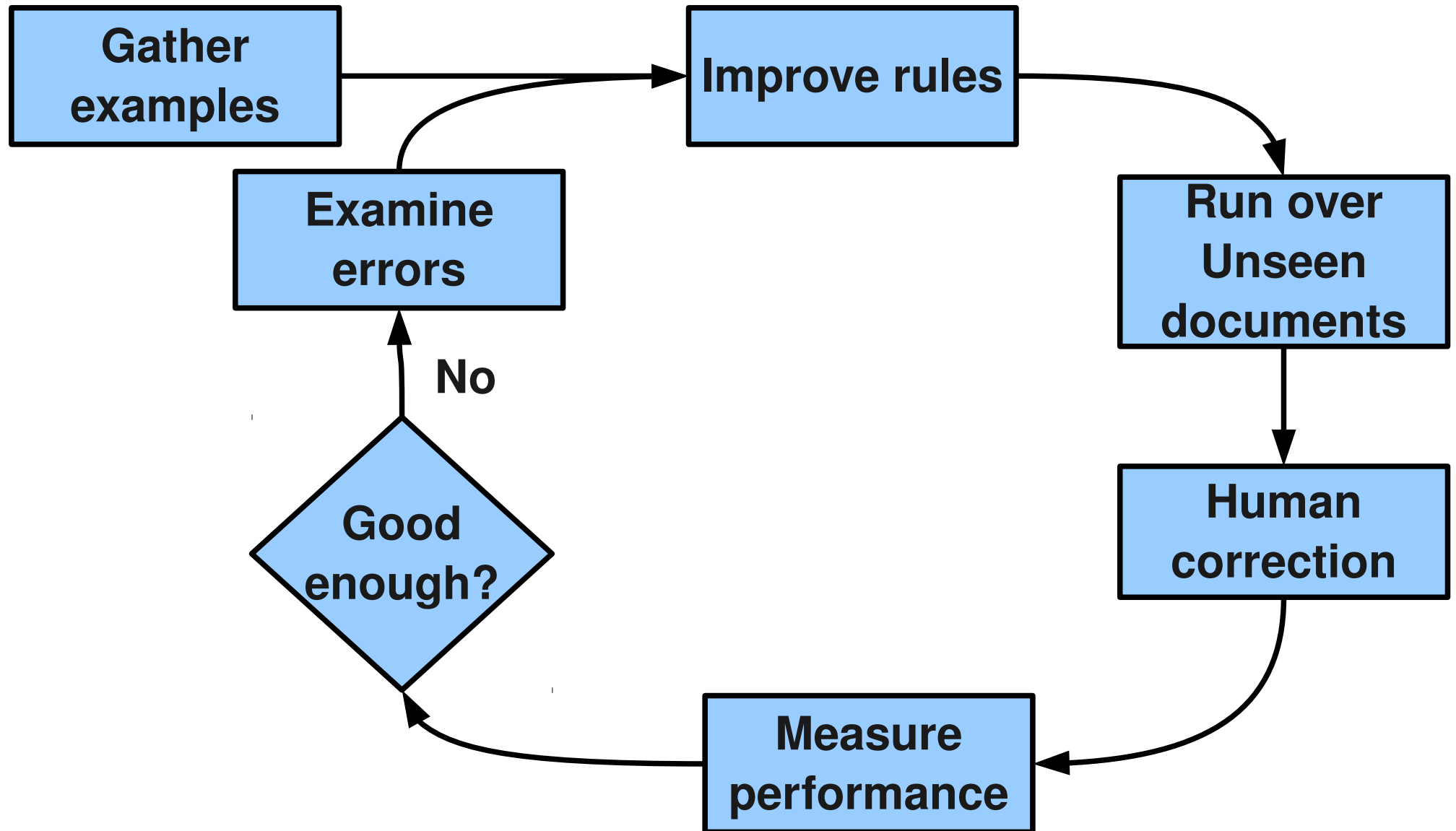
Missing: 2 Lenient: 0.6667 1.00 0.80

False Positives: 0 Average: 0.5833 0.875 0.70

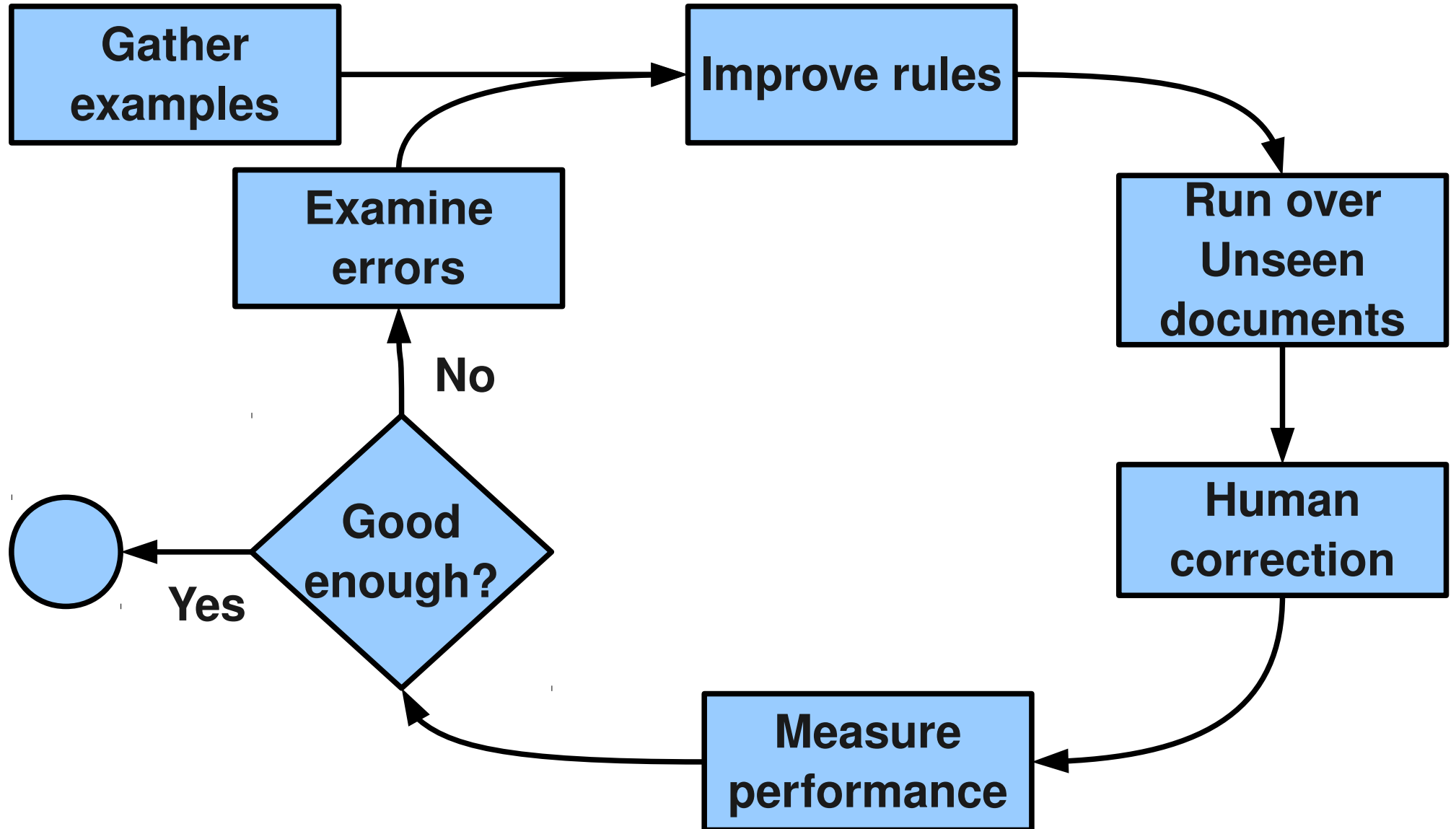
A process



A process



A process



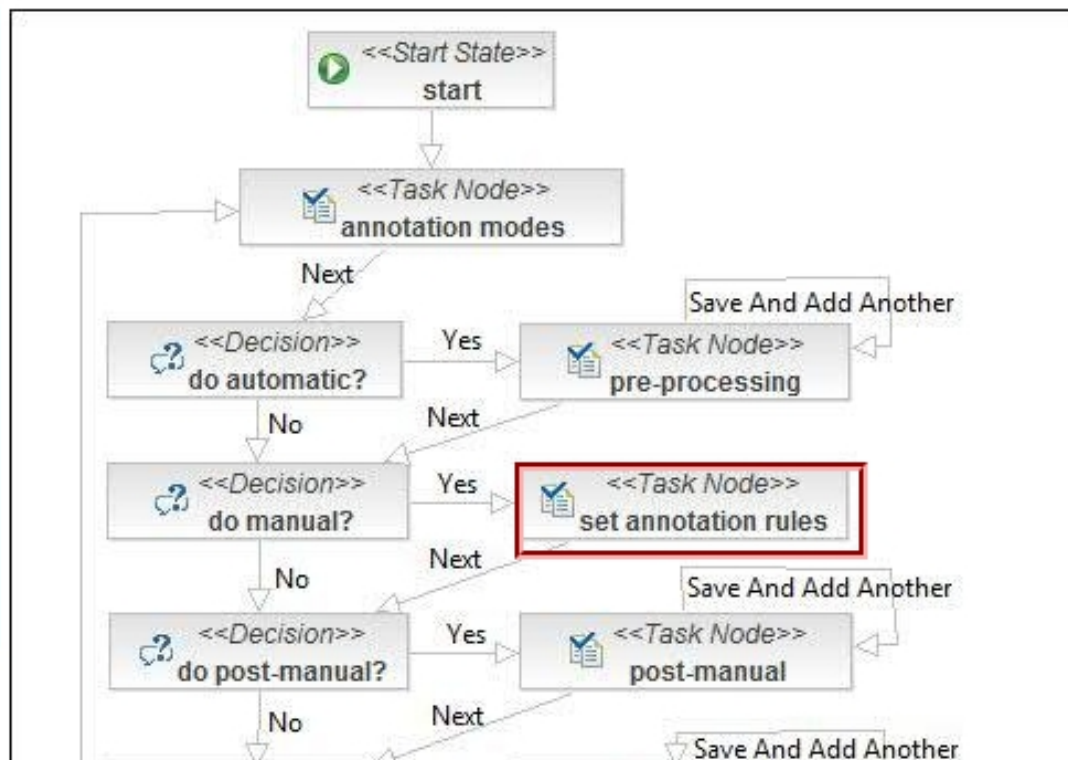
Supporting the process

- We need to train and implement the process
- We need tools to support this process
- Quality Assurance Tools
- Workflow: GATE Teamware
- Annotation pattern search: Mimir
- Coupling search and annotation: Khresmoi
 - Develop a pattern and run it over text
 - Human correction
 - Feedback

GATE Teamware: defining workflows

Annotators per Document:	2
Cancel Task Allowed:	<input checked="" type="checkbox"/>
Anonymous Annotation:	<input checked="" type="checkbox"/>
Annotation Schemas:*	<ul style="list-style-type: none"> OrganizationSchema.xml PersonSchema.xml AddressSchema.xml Non-FunctionalRequirements ConcordiaDataGroupsSc ConcordiaRequirementsS
Add Schema:	Add Schema
Pre-Manual Service:*	None

Next Help Quit



GATE Teamware: managing projects

Project Name * <input type="text" value="Manual WF project"/>	Managers <input type="text" value="kalina"/>	Curators <input type="text" value="ac4sa59"/> <input type="text" value="agaton"/> <input type="text" value="ishrar-cur-man"/> <input type="text" value="kalina"/> <input type="text" value="kalina-curator"/> <input type="text" value="matthew-cur-man"/> <input type="text" value="milan"/>	Annotators <input type="text" value="ac4sa59"/> <input type="text" value="adam"/> <input type="text" value="agaton"/> <input type="text" value="angus"/> <input type="text" value="ayrin-ann"/> <input type="text" value="danica"/> <input type="text" value="diana"/> <input type="text" value="hamish"/> <input type="text" value="ian"/> <input type="text" value="ishrar-ann"/>
---	--	---	--

[Add Corpus](#)

Save & Start

My Workflow Templates

GATE Teamware: monitoring projects

Process Monitoring: Annotation Status

Detailed View

Back to Project

Status	#
Annotated	64
Canceled	1
Failed	0
In Progress	1
Not Started	7

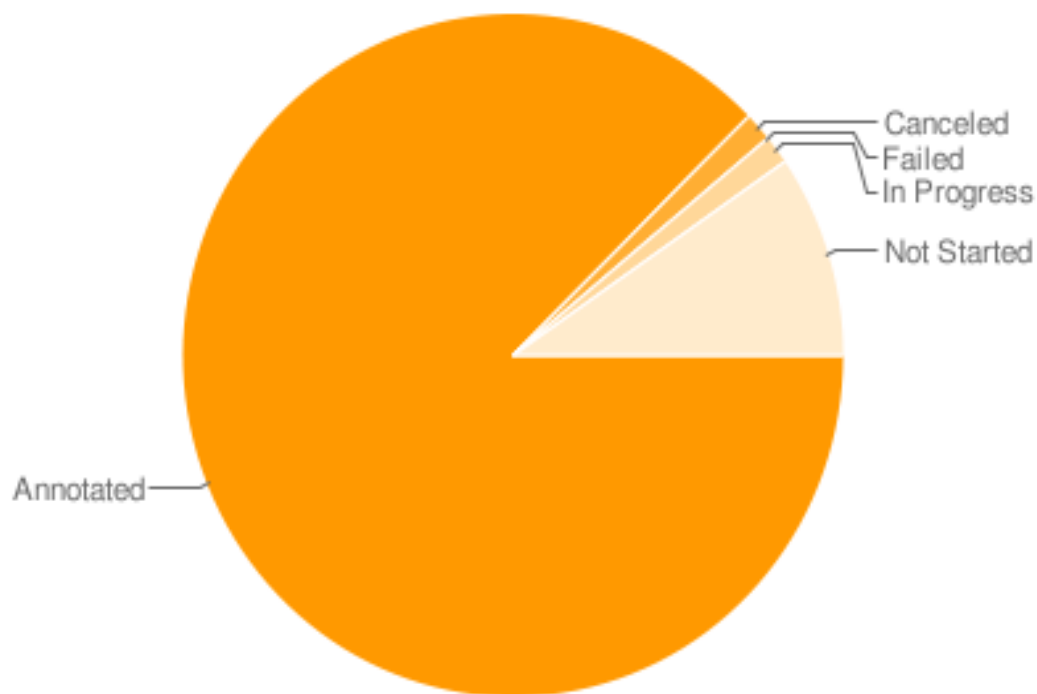
Average Execution Time

627.609375

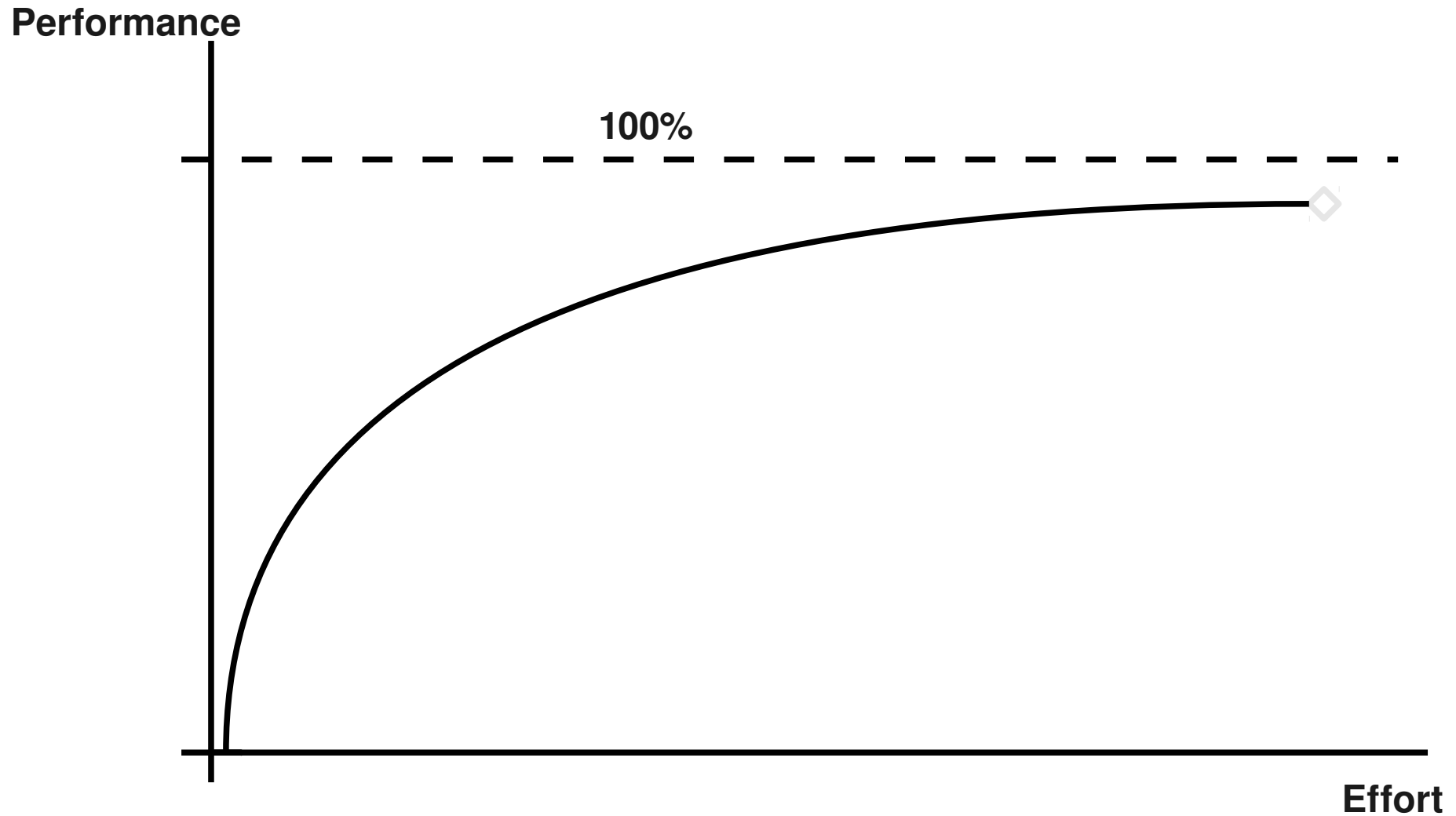
Detailed View

Back to Project

Annotation Status Chart



When is it good enough?



When is it good enough?

- Like disease
 - A small number of very common ones
 - Lots of rare ones
- For a straightforward use case,
 - 3 or 4 iterations
 - plateau at around 90%

Implications

- Ad-hoc annotation and search is as important an approach as generic annotation
- We need tools and processes to support this style of annotation
- An agile annotation process involves users, helps us to elicit their requirements, and reduces the cost of annotation