

Crowdsourcing Approach in Statistical Machine Translation

Project Title

Development of Multilingual Resources and Technologies for Indian Languages

(a collaborative project between IIT Bombay and Xerox)

Presented By:- Rajen Chatterjee (Research Engineer, IIT Bombay)

Team Members:- Abhijit Mishra (Research Scholar, IIT Bombay)

Anoop Kunchukuttan (Research Scholar, IIT Bombay)

Outline

- What is Crowdsourcing?
- Role of Crowdsourcing in SMT
- How to manage crowd?
- Basic terminologies
- Need for automation?
- Crowdsourcing Platform Architecture
- Crowdsourcing Engine
- Workflow
- Different Stages
- Pipelining and Quality Control
- A working example

What is Crowdsourcing?

- **Crowdsourcing** is a process that involves outsourcing tasks to a distributed group of people.
 - Example: Recaptcha used in Gmail account creation.
- This process can occur both online and offline.
- The difference between crowdsourcing and ordinary outsourcing is that a task or problem is outsourced to an undefined public rather than a specific body.

Role of crowdsourcing in SMT

- Crowdsourcing approach is being used for development of Parallel Corpora.
- Not just development of parallel Corpora but development of qualitative parallel Corpora.
- Ensuring quality is the most challenging task in Crowdsourcing.
- Crowd Tasks:
 - Phrase Translation
 - Phrase Validation
 - Phrase Re-Combination

How to manage crowd?

- Third party service provider like Amazon Mechanical Turk, CrowdFlower.
- Amazon Mechanical Turk is a web service that provides an on-demand, scalable, human workforce to complete jobs that humans can do better than computers.
- Sample Business Cases:
 - Translation
 - Image Tagging
 - Data Verification
 - Information Gathering
- Service Highlights:
 - On-Demand Workforce
 - Quality Management Tools
 - Lower Cost Structure

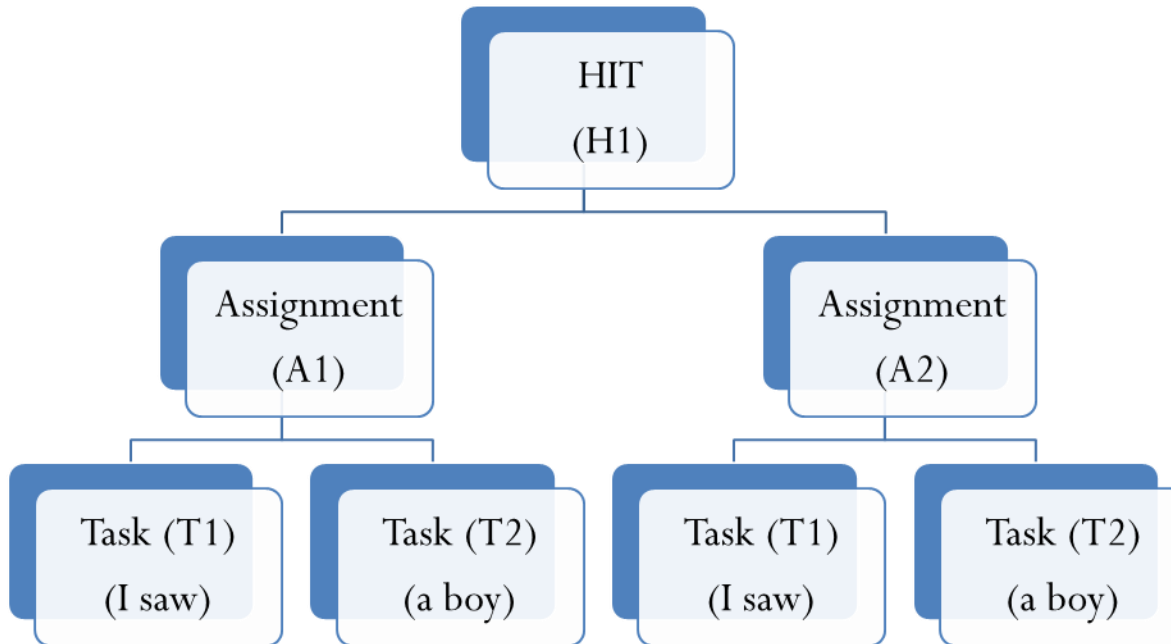
Example..

- Consider the following small sentence

I saw a boy.

- Our objective is to get the Hindi translation of the whole sentence.
- In first phase, will divide it into two phrases viz. “I saw” and “a boy” and will obtain translations of the corresponding phrases from the crowd.
- Second phase will ensure whether the quality of the target phrase is good or not.
- If the quality is good, we’ll recombine the phrases in the final phase to get the whole target sentence.

Basic Terminologies



- Task:
 - Task is the atomic unit of work.
- Assignment:
 - Several task are combined into one Assignment.
- HIT (Human Intelligence Task):
 - Several Assignment are combined into one HIT.

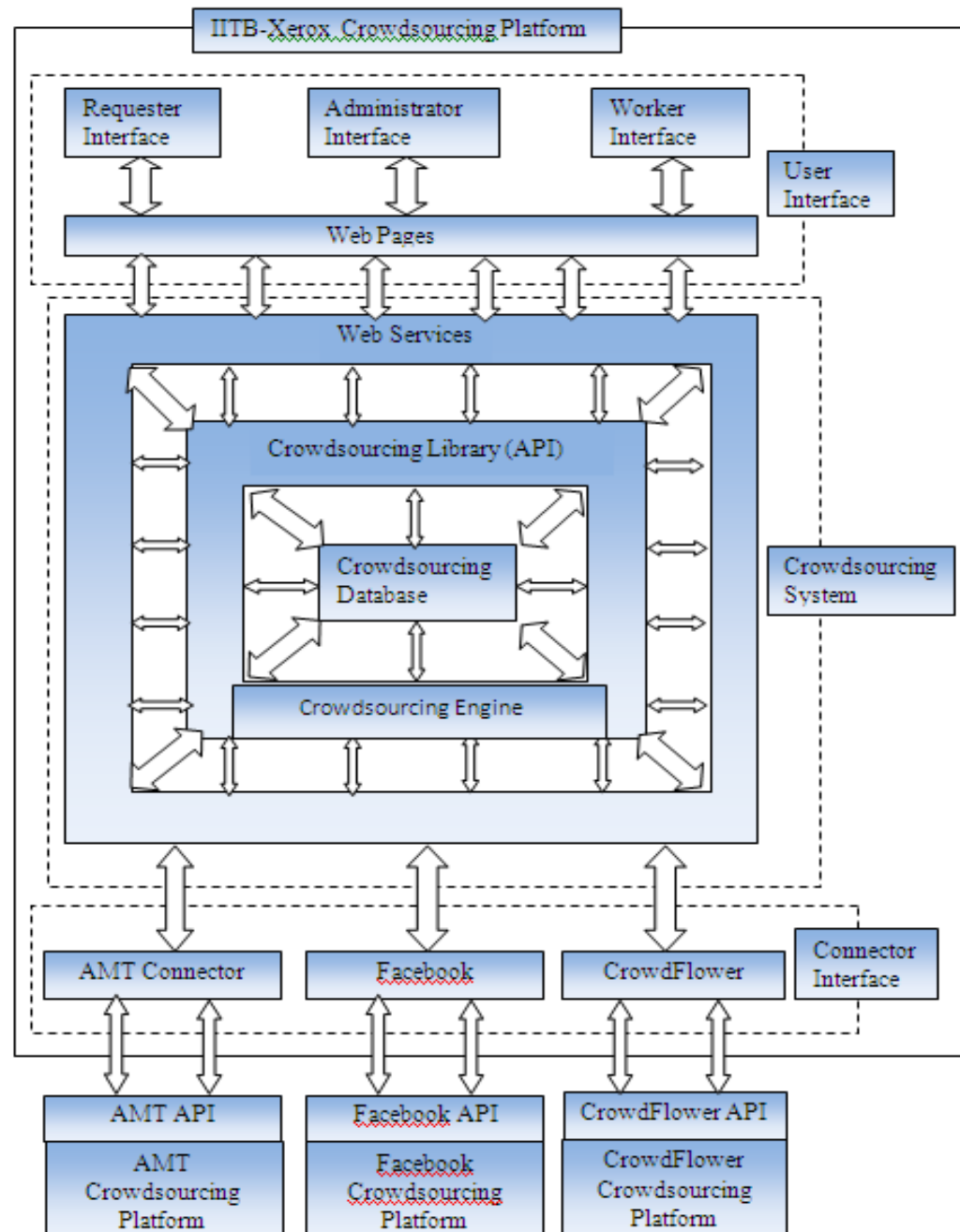
Need for Automation?

- Manually managing Crowdsourcing task is very tedious specially in case of Translation.
- So there is a need to automate the whole crowd translation process.
- Automation will help to get the following things done:
 - Get the Translation of each source phrase
 - Validate the translated phrase.
 - Re-combine the translated phrase into sentence.
- This brings up the motivation for developing Crowdsourcing Platform.

Crowdsourcing Platform

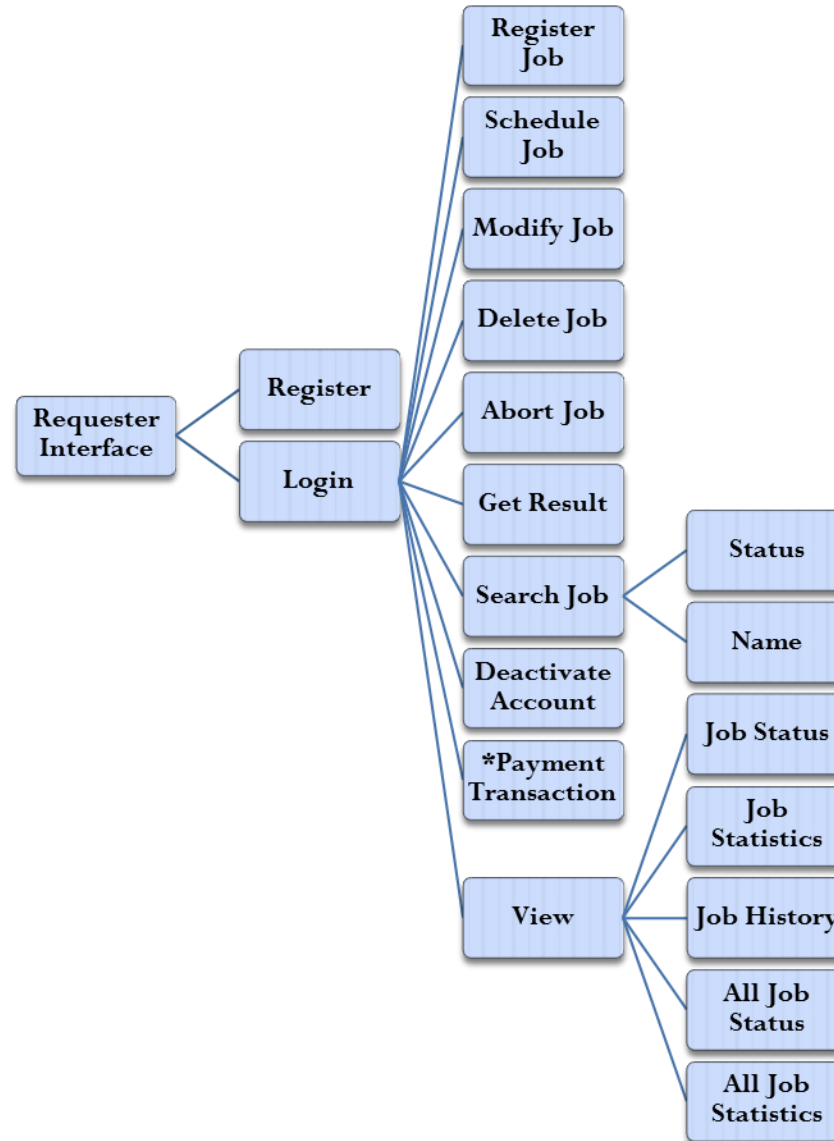
- Goal of this platform is to develop a generic task independent platform but focusing more on Translation task.
- From translation point of view our goal is to develop language independent and domain independent parallel corpora.

ARCHITECTURE



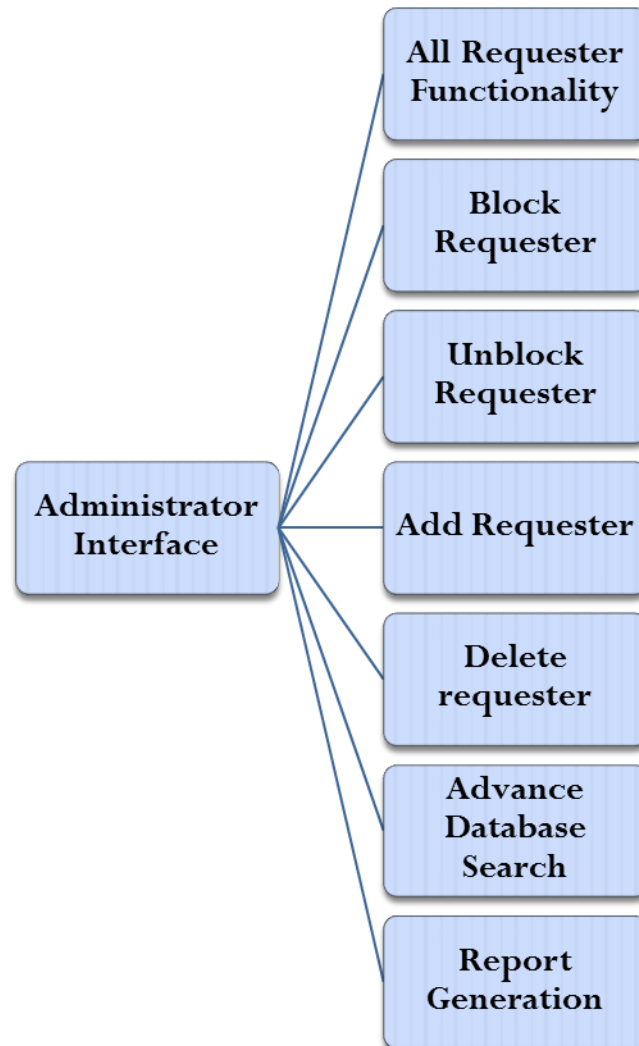
Requester

- Requester is the one who requests for getting a translation of a given document.



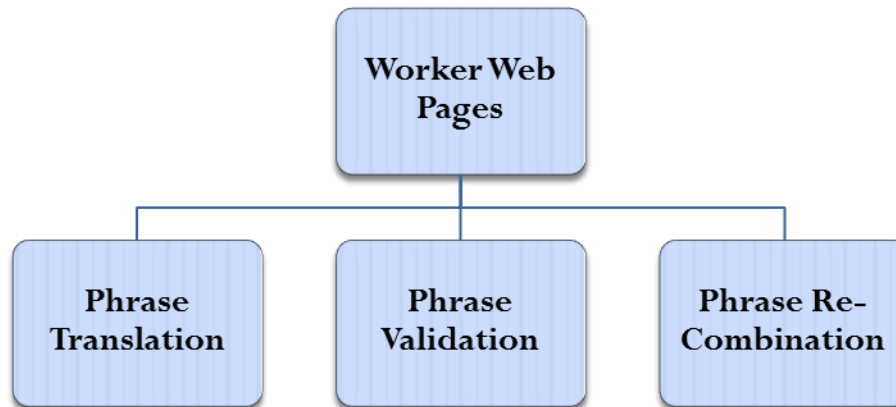
Administrator

- Administrator will have all rights to perform all tasks on behalf of a requester



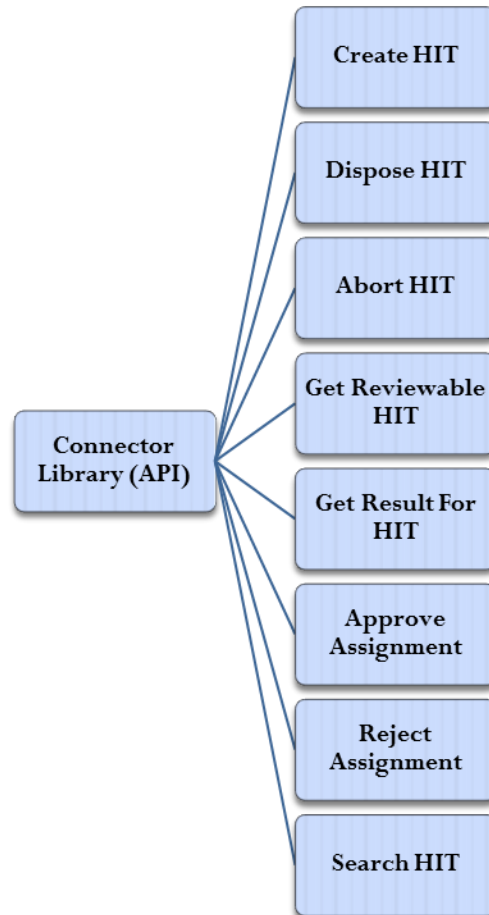
Worker

- Workers are the member of crowd itself who are actually involved in solving the task. Workers will be connected to Crowdsourcing worker web page indirectly via other Crowdsourcing platform like AMT.
- Workers will be provided with any one of the following web pages for translation job.
 - Phrase Translation
 - Phrase Validation
 - Phrase Re-Combination



Platform Connector

- Platform connector behaves as an interface between IITB-Xerox Crowdsourcing system and other Crowdsourcing platform.
- Each connector will have its own API functionality which will interact with connector Crowdsourcing platform.



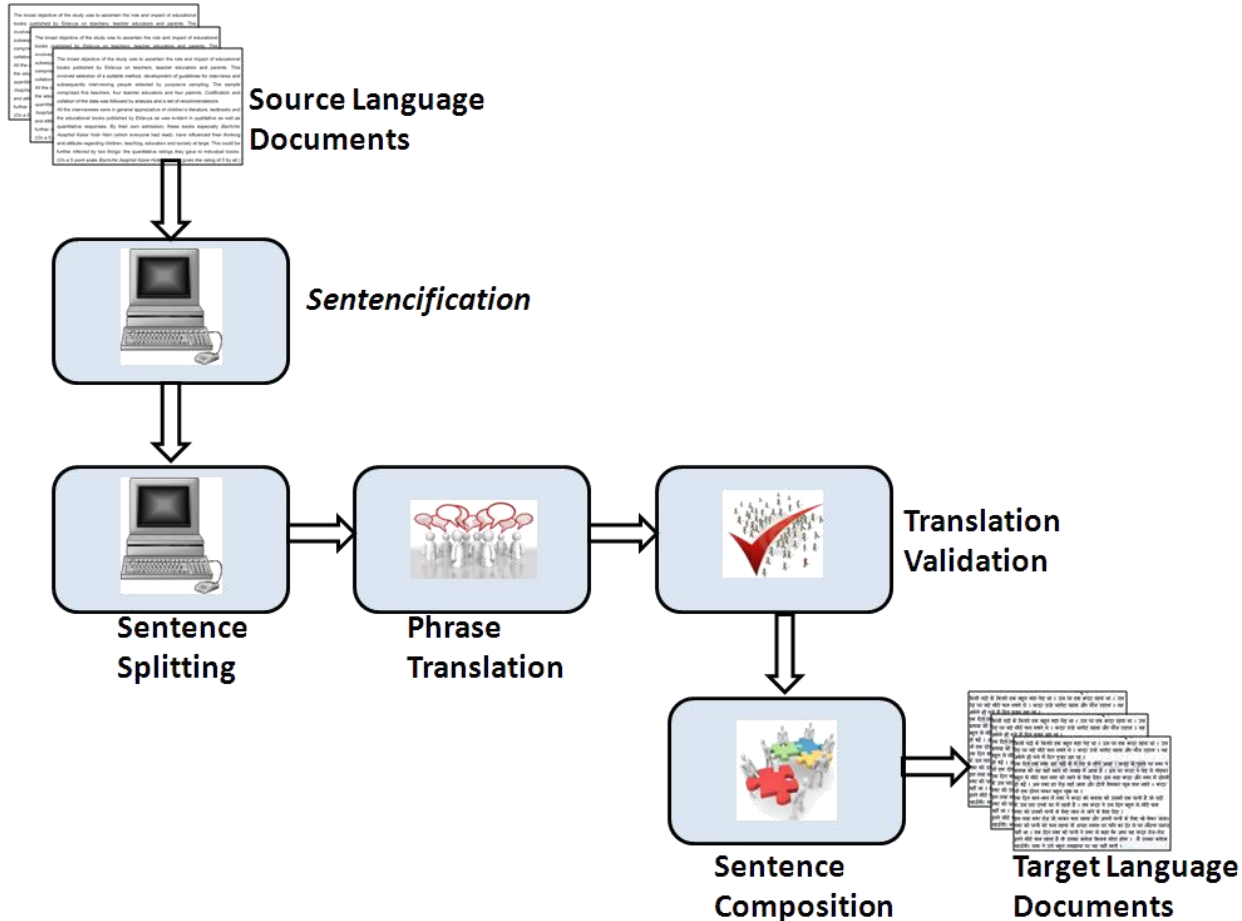
Platform Connector

- Following are the role of a connector:
 - Pulling And Pushing Mechanism:
 - Pulling data from IITB-Xerox Crowdsourcing system and pushing into AMT.
 - Pulling result from AMT and pushing into IITB-Xerox Crowdsourcing system.
 - Monitoring:
 - Periodically monitoring IITB-Xerox Crowdsourcing system to check any status change of a HIT
 - Periodically monitoring AMT Crowdsourcing platform to check for the result of HIT.

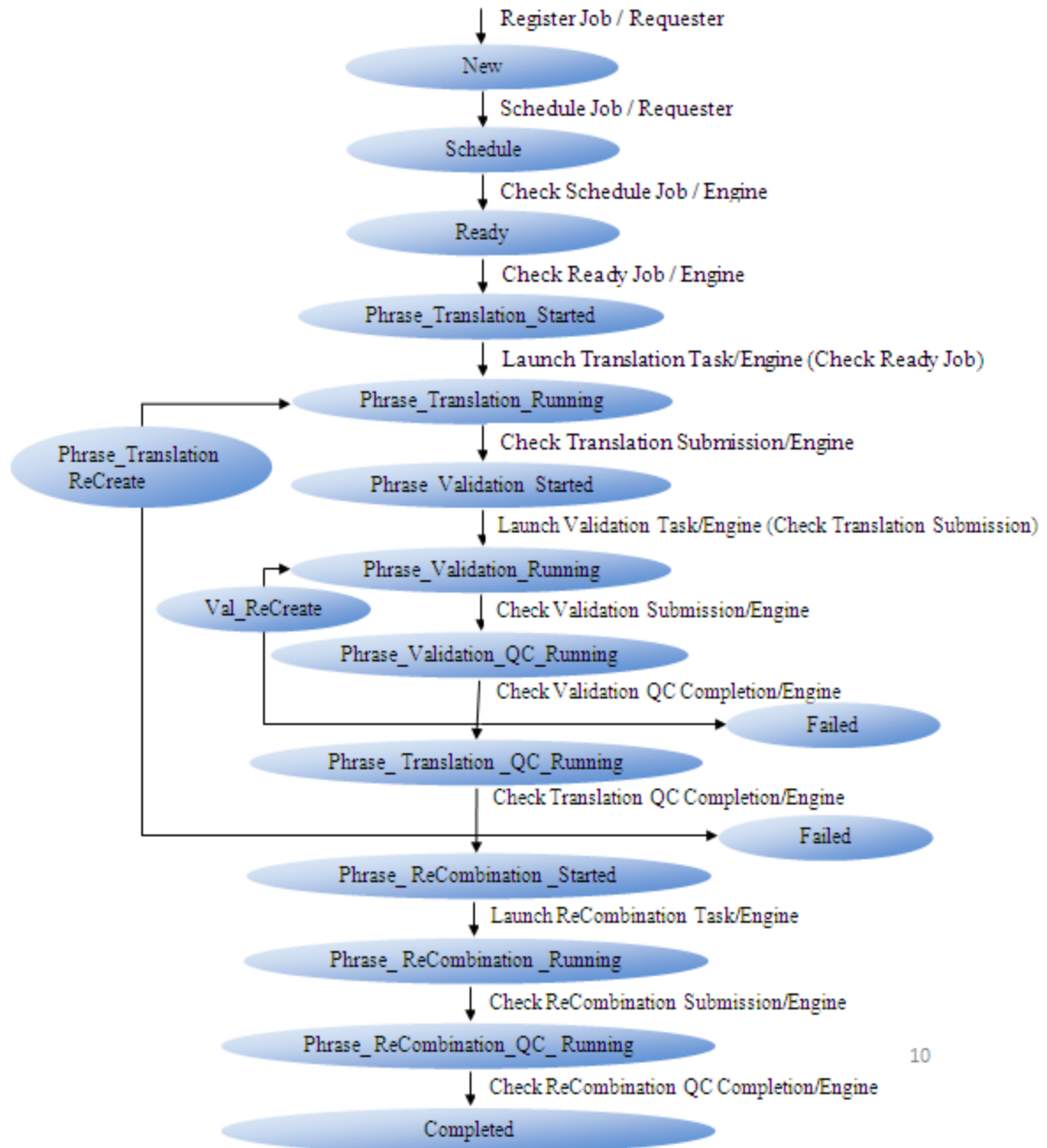
Crowdsourcing Engine

- Crowdsourcing Engine is the heart and soul of the Crowdsourcing Platform.
- Following are the responsibility of the engine:
 - Monitoring entire life cycle of a job.
 - Managing the Pipelining architecture.
 - Ensuring Quality Control Mechanism.
 - Update the HIT status whenever Job status changes.

Workflow (A Naïve view)

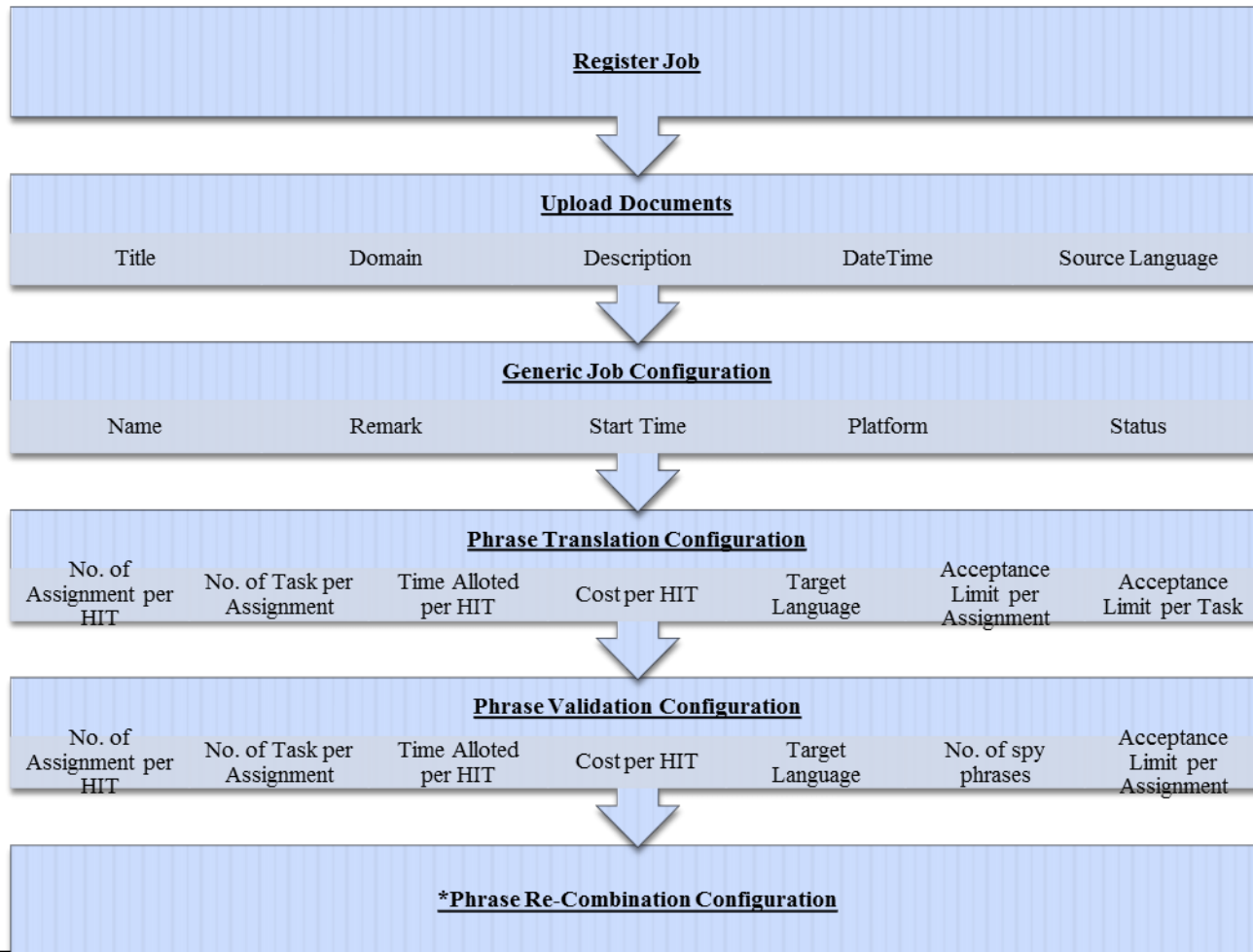


W
O
R
K
F
L
O
W

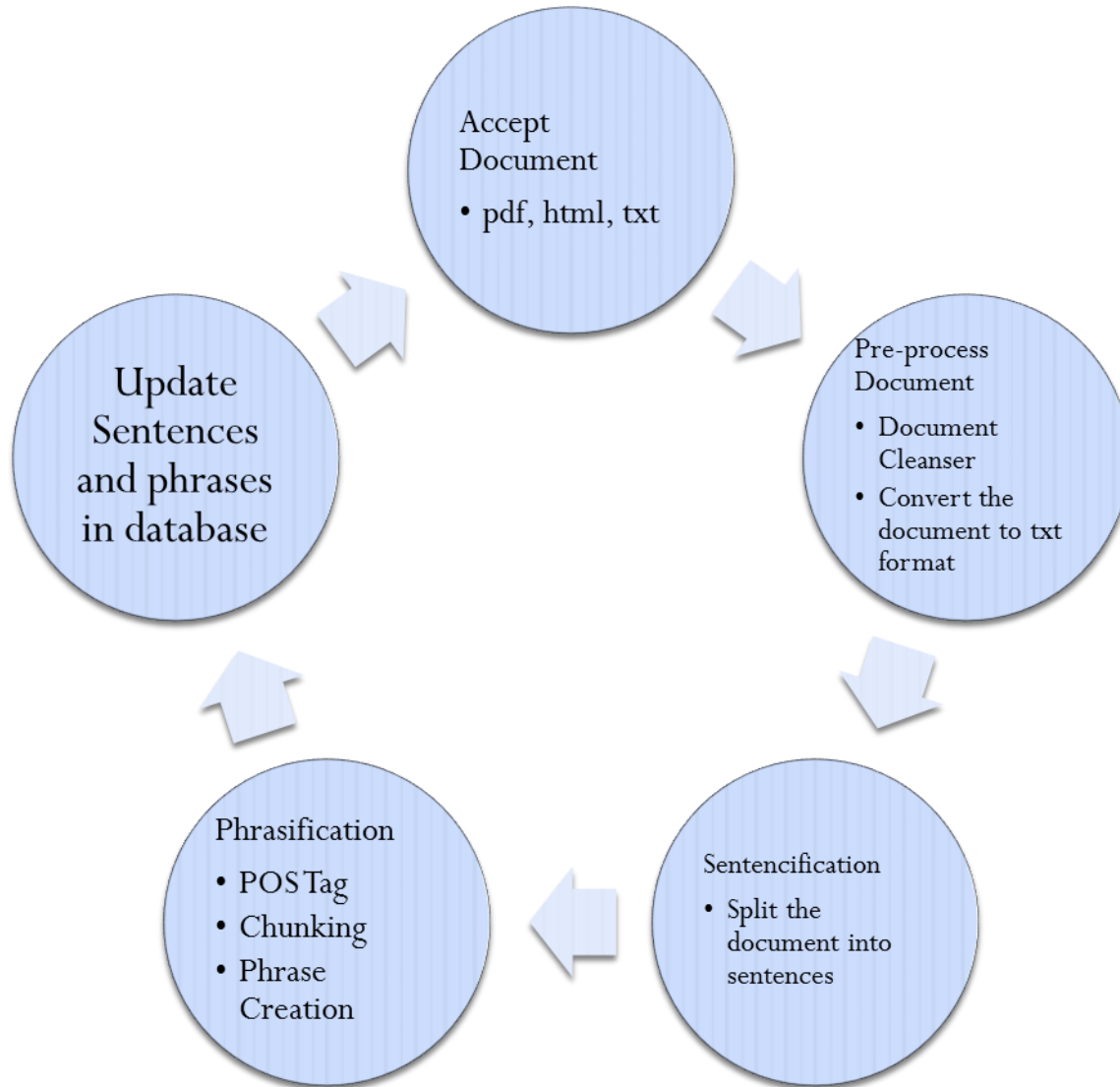


Register job

- Register job will only create a job for execution i.e. it will register all job configuration parameters in the database and upload the documents for translation



Document processor



Pipelining & Quality Control

- Pipelining ensure that different stages in entire translation process runs smoothly.
- Quality Control ensures that we get the translated phrase of reasonable quality.

Phrase Translation

- In Phrase Translation, given a source phrase the worker is required to give the translated phrase.

The screenshot shows the 'English to Hindi Phrase Translation' interface. At the top, a header bar displays 'English to Hindi Phrase Translation', 'Requester: Translation Test', 'Qualifications Required: None', 'Reward: \$0.05 per HIT', 'HITs Available: 228', and 'Duration: 60 minutes'. The main content area is titled 'English to Hindi Phrase Translation' and includes 'PREVIOUS' and 'NEXT' buttons, and a progress indicator '2 of 5 phrases'. The source text 'are available for such purchases and payments' is on the left. A text input field labeled 'उपलब्ध' is on the right. Below, 'Suggested Translations' are listed: 'available : उपलब्ध', 'purchase : खरीद क्रय करना खरीदना क्रय', 'and : और एंड', 'for : के लिए', and 'are : हैं क्या'. Callouts point to 'Use of images' (referring to the source text), 'Transliterated input' (referring to the input field), 'Uncluttered display' (referring to the clean layout), and 'Vocabulary support' (referring to the suggested translations).

English to Hindi Phrase Translation

Requester: Translation Test
Qualifications Required: None
Reward: \$0.05 per HIT
HITs Available: 228
Duration: 60 minutes

Instructions (Show)

Use of images

English to Hindi Phrase Translation

PREVIOUS NEXT

2 of 5 phrases

Transliterated input

are available for such purchases and payments

उपलब्ध

Suggested Translations

available : उपलब्ध

purchase : खरीद क्रय करना खरीदना क्रय

and : और एंड

for : के लिए

are : हैं क्या

Uncluttered display

Vocabulary support

Phrase Validation

- In Phrase Validation, for a given source and target phrase the worker is required to assign a rating to the translation on a scale of 5 (0 means worst and 5 means best)



An Open Innovation Initiative by Xerox Innovation Group and IIT Bombay
Developing Linguistic Resources for Indian Languages
Your chance to Beat the Machine!



Instructions [\(Hide\)](#)

- To begin the task you need to first click on the **"Accept HIT"** button at the top/bottom of the page.
- Kindly verify whether the phrasal translations given below are correct or not.
- Please select **"No"** in case you are not sure.

English Phrase	Hindi Translation	Is the translation correct?
An additional copy	और एक कापी	<input type="radio"/> Yes <input checked="" type="radio"/> No
be sent to the appellant	अपीलकर्ता को भेजा गया	<input checked="" type="radio"/> Yes <input type="radio"/> No
through the Superintendent , Central Jail , Tihar .	तिहाड़ जेल के प्रबंधक के द्वारा	<input checked="" type="radio"/> Yes <input type="radio"/> No
of the moment .	पल की	<input checked="" type="radio"/> Yes <input type="radio"/> No
about the authenticity of this letter .	इस पत्र की प्रामाणिकता के बारे में ।	<input checked="" type="radio"/> Yes <input type="radio"/> No
against the petitioner .	आवेदक के विरुद्ध में	<input checked="" type="radio"/> Yes <input type="radio"/> No
than 30 years old .	30 से अधिक वर्षों पुरानी	<input type="radio"/> Yes <input checked="" type="radio"/> No
by the Arbitrator	पंच द्वारा	<input checked="" type="radio"/> Yes <input type="radio"/> No
of claimants in this case	मुकदमे में दावेदारों का	<input checked="" type="radio"/> Yes <input type="radio"/> No
Whether reporters of the local news papers	अगर स्थानीय अखबारों के संवाददाताओं को	<input checked="" type="radio"/> Yes <input type="radio"/> No

Submit

Phrase Re-combination

- In Phrase Re-combination, given a list of target phrases the worker is required to align the phrases and do necessary editing to obtain a valid target sentence.

Build Hindi Sentences from Phrases

Instructions ([Show](#))

An additional copy be sent to the appellant through the Superintendent , Central Jail , Tihar .

Workspace (Reorder the phrases by dragging them down onto the gray boxes)

अपीलकर्ता को भेजा जा

Hindi phrases re-ordered

एक अतिरिक्त प्रतिलिपि तिहाड़ जेल के प्रबंधक के द्वारा

Edit reordered sentence

Reset

Final Hindi Sentence

Submit

Quality Control for Validation

- Spam Detector/ Noise Removal:
 - Each Validation Assignment will contain few spy phrases.
 - If a worker fails to successfully validate this spy phrases then the worker will be consider as spam and the Assignment will be rejected
 - The rejected Assignment will be resend to the crowd by the pipelining module.
 - After certain number of iteration if we are not getting valid assignment then respective measures will be taken.
 - If we get valid assignment then the rating will be used for analyzing quality of phrase translation.

Quality Control for Translation

- Based on the validation rating for same pair of source and target phrase an average score will be computed and assigned to that target phrase.
- The target Phrase which will have the highest score for a particular source phrase will be considered as the winner.
- If the winner target phrase is above a threshold value then it will be preserved else that source phrase will be resend to crowd for translation.
- After certain number of iteration if we are unable to get a good translation then respective measure will be taken.

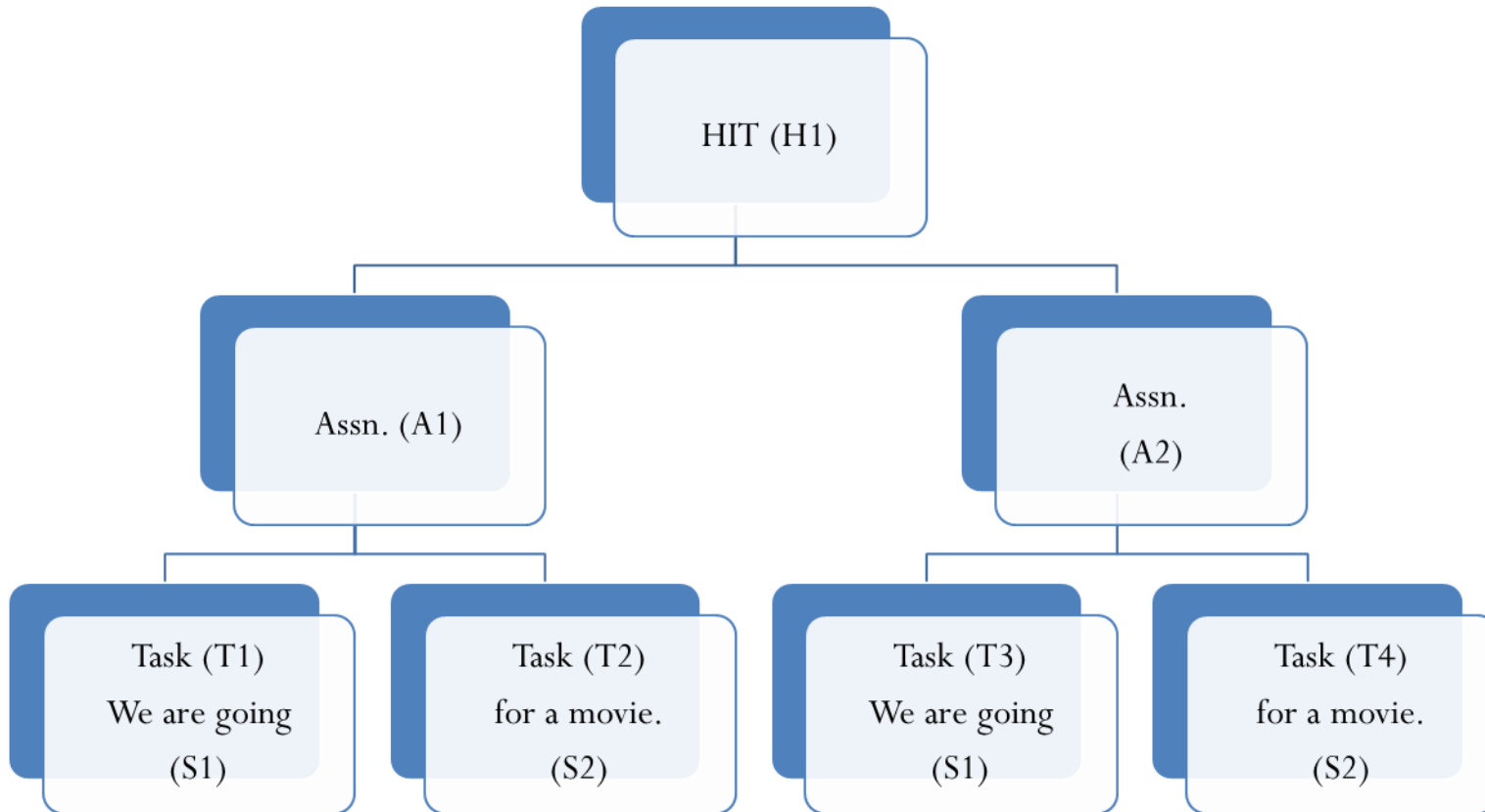
Quality Control for Re-combination

- Quality Control for Re-combination begins once the phrase re-combination phase is over.
- It checks for the proper alignment of phrases and whether the sentence is complete on its whole.
- Linguist may be introduced in this stage to get an authorized confirmation on the translated sentences.

Working example

Phrase Translation

- ⊙ Sentence = We are going for a movie.
- ⊙ S1 = We are going, S2 = for a movie.

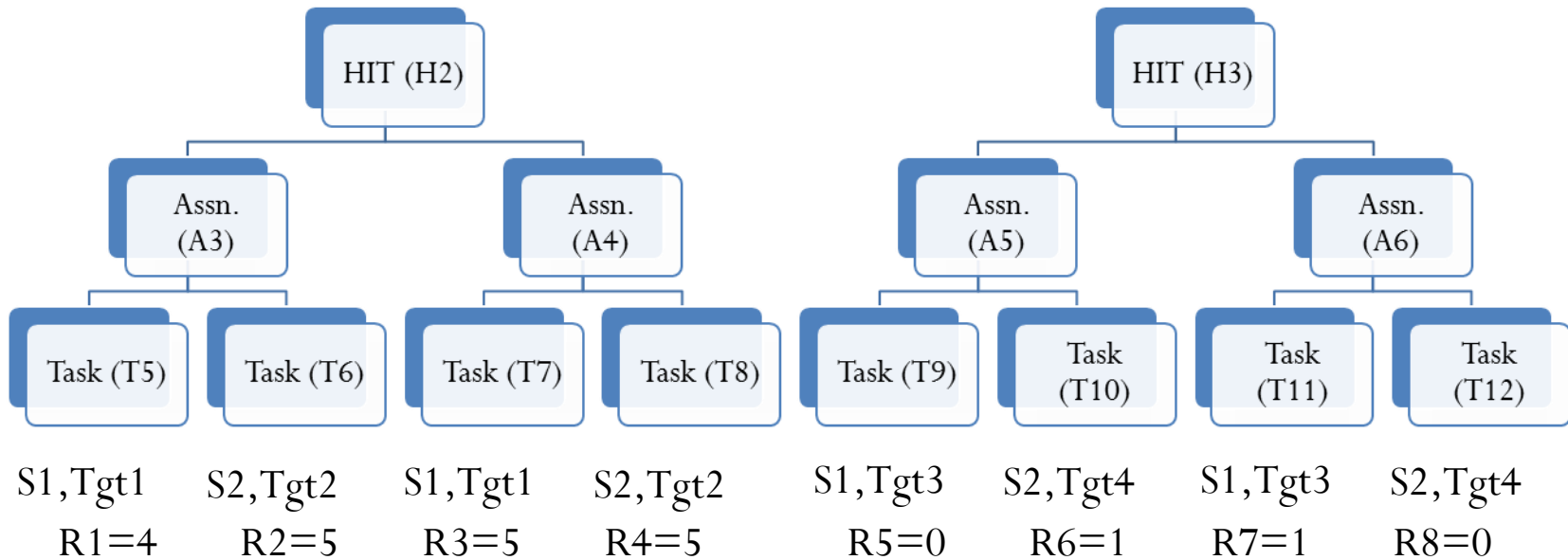


Tgt1 = हम जा रहे हैं Tgt2 = मूवी के लिए Tgt3 = वे अरे गोइंग Tgt4 = फॉर अ मूवी

Phrase Validation

S1=We are going S2=for a movie.

Tgt1= हम जा रहे हैं Tgt2= मूवी के लिए Tgt3= वे अरे गोइंग Tgt4= फॉर अ मूवी



$$S1, Tgt1 = (R1+R3)=(4+5)/2=4.5$$

$$S1, Tgt3 = 0.5$$

$$S2, Tgt2 = 5$$

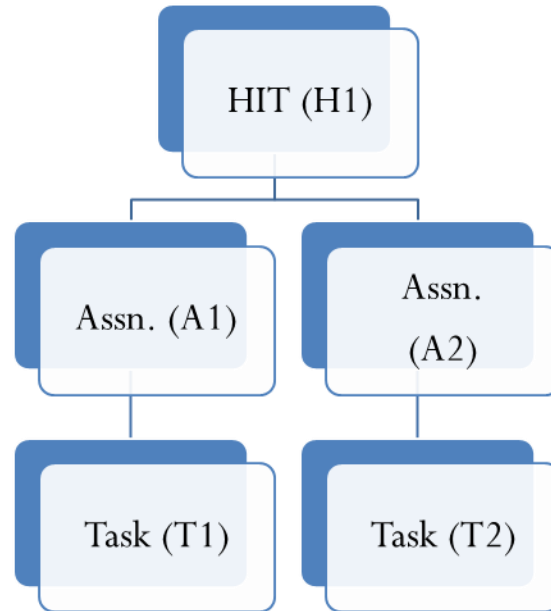
$$S2, Tgt4 = 0.5$$

Phrase Re-combination

Source Sentence=We are going for a movie

Tgt1= हम जा रहे हैं

Tgt2= मूवी के लिए



हम जा रहे हैं मूवी के लिए

हम मूवी देखने जा रहे हैं

<http://www.cse.iitb.ac.in/~abhijitmishra/reorder.html>

Time for a Demonstration

Demonstration Steps:-

