

Acquisition of Language Symbols

Pankaj Prateek (pratikk@)
Mentor: Dr. Amitabha Mukerjee (amit@)
Dept. of CSE, IIT Kanpur

Previous Work

The baby designer model learns patterns in an apprenticeship situation. When presented with a set of functional constraints and variable set governing them, it explores the design space, using domain-general learning algorithms to discover patterns in the better performing designs. These patterns get transformed to chunks in case they occur frequently. In this process, knowledge of language and labels for such concepts and patterns is not required.

On exposure to language, these implicit associations get transformed to rules in the symbolic space, thus incorporating labels.

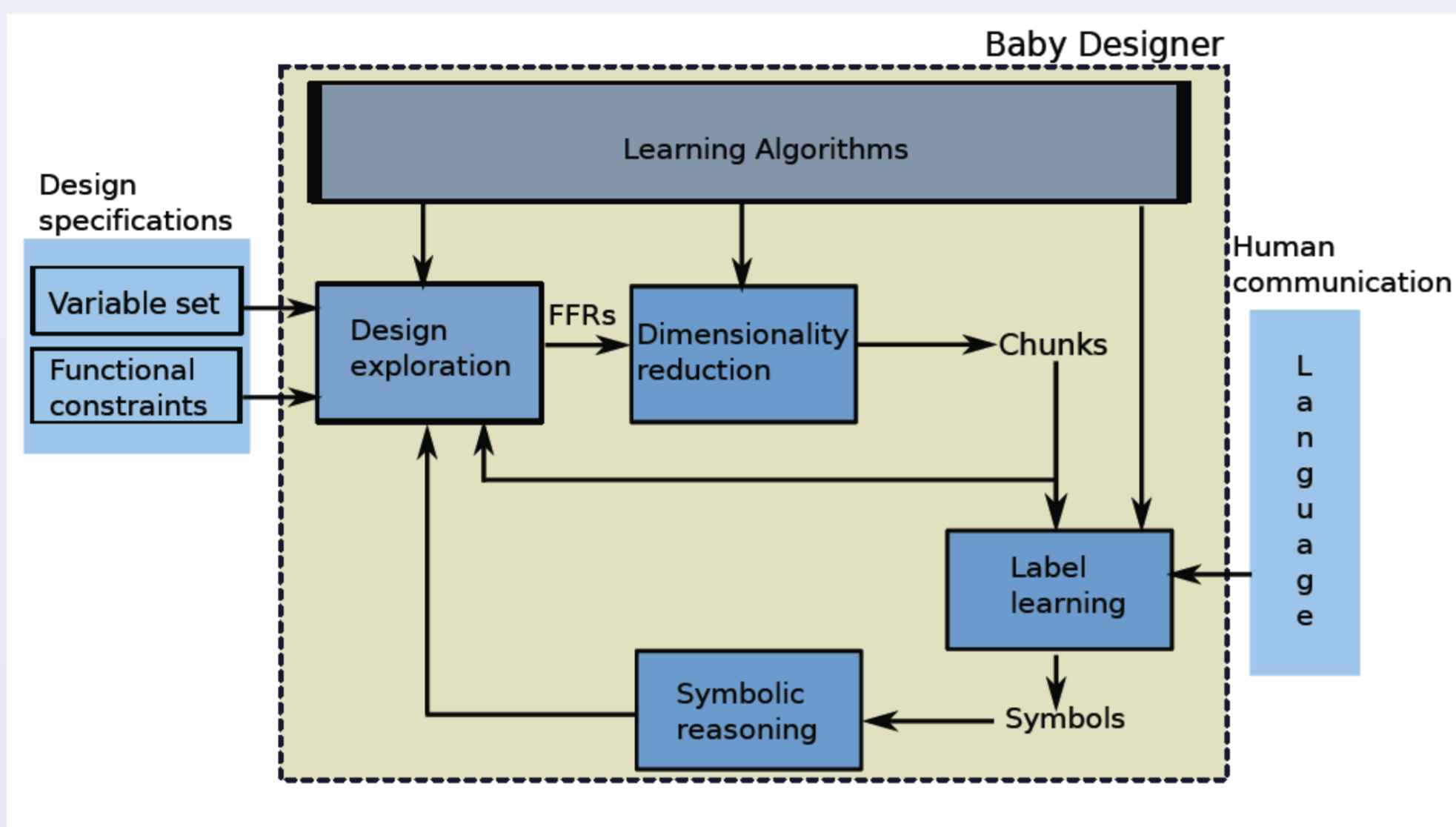


Fig 1: Architecture of the Baby Designer

Due to difference in experience and language, the symbols acquired by different agents differ. In their paper[1], Dr. Mukerjee and Madan Dabberu have considered how agents map the chunks (low-dimensional characterizations) to language based on human commentary produced in the same context. This was explored by learning labels in English and Telugu in the simple domain of tight and loose fits using the peg-in-hole assembly. This project aims to extend their work to Hindi.

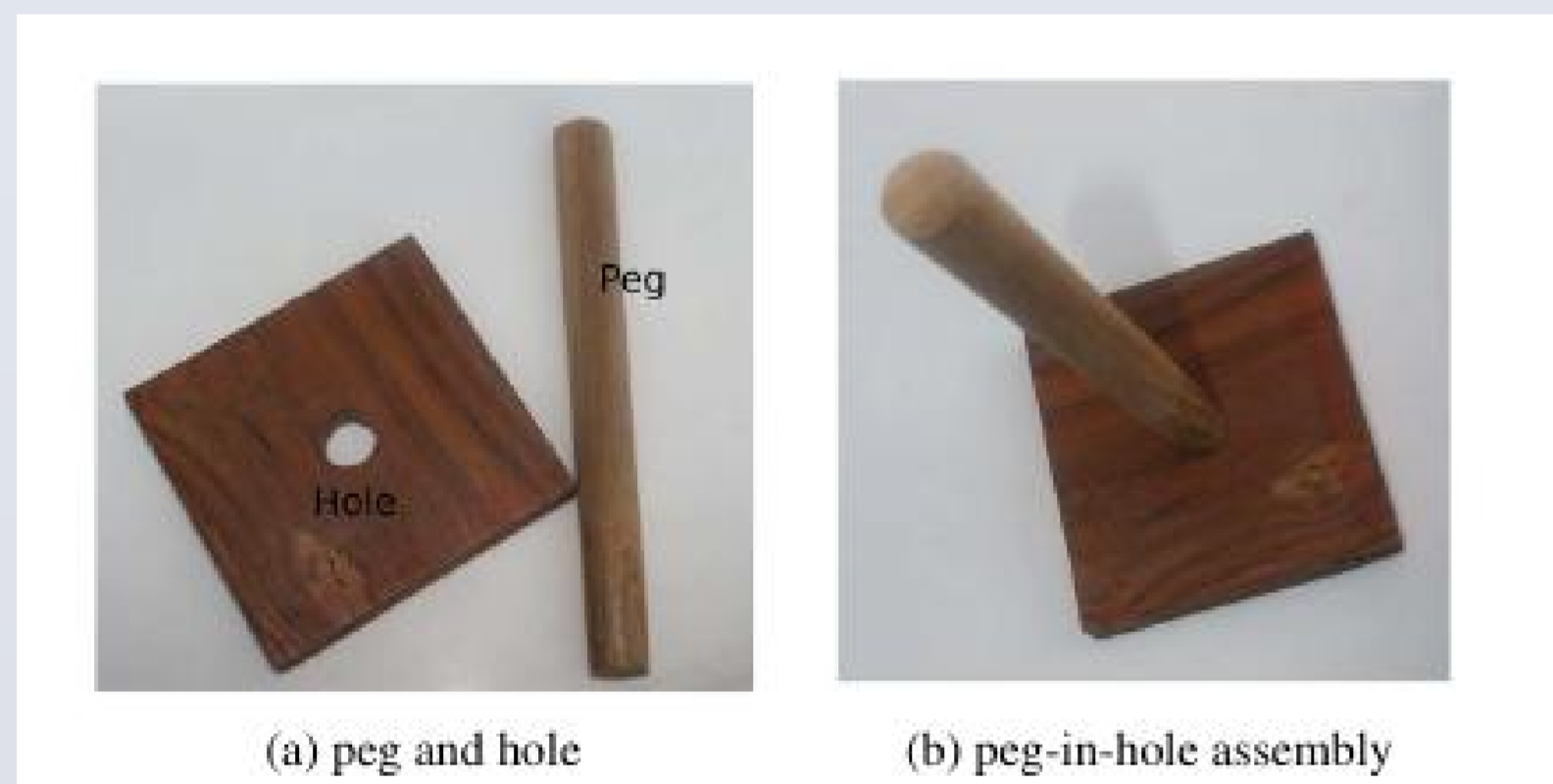


Fig 2: Peg-in-hole Assembly

Apparatus

A wooden block with 5 Holes (A (tapering), B, C, D, E) and 5 aluminium/steel pegs (1,2,3,4,5) of diameters (9mm, 12.25mm, 15.75mm, 19mm, 25mm).

Loose-fit situations - A:1, C:5, D:3

Tight-fit situations - A:2, B:4, E:3

Methodology

Experiment 1:State

This experiment focused at collecting spoken Hindi data in situations when the peg is already inserted in the hole. The pre-assembled setup was kept on a table and the participants were asked not to lift the assembly from the table.

Each participant was provided with the following instruction:

“यह एक छेद है और यह एक पैग है. यह पैग पहले से इस छेद में डाली हुई है. आपको इन दोनों के बीच हो रहे संपर्क का हिन्दी में वर्णन करना है. यह ध्यान रखिये कि यह ब्लॉक इस टेबल से उठे नहीं और जितना हो सके हिन्दी के अलावा किसी और भाषा का प्रयोग ना करिये.”

Experiment 2: Action

This experiment focused at collecting spoken Hindi data in situations where the participant is asked to actively insert the peg in the hole. There was no restriction as in the previous experiment.

Each participant was provided with the following instruction:

“यह एक छेद है और यह एक पैग है. आपको यह पैग इस छेद में डालनी है और इन दोनों के बीच हो रहे संपर्क का हिन्दी में वर्णन करना है. जितना हो सके हिन्दी के अलावा किसी और भाषा का प्रयोग ना करिये.”

Associating Linguistic Labels

The association of a word w_i with a concept C_j can be measured using conditional probability. But the direction of association could be either $p(C/w)$ or $p(w/C)$.

Since only two concepts are involved:

$$\frac{p(C_T/w)}{p(C_L/w)} = \frac{p\left(\frac{w}{C_T}\right) \cdot p(C_T)}{p\left(\frac{w}{C_L}\right) \cdot p(C_L)}$$

The number of instances of C_L and C_T are almost same in the set, so the direction of association doesn't matter.

For strongest association with C_T , compute $\max_i \left\{ \frac{p\left(\frac{w_i}{C_T}\right) \cdot p(C_T)}{p\left(\frac{w_i}{C_L}\right) \cdot p(C_L)} \right\}$

Results

Words like “tight” and “loose” are readily associated with [tight] and [loose] fits, respectively, based on uninformed word association from Hindi Language. Thus the words have top associations with the concepts.

Unlike Telugu[1], Hindi is richly inflected, it incorporates foreign words very easily into the language. This can be seen from the fact that the Hindi native correspondences for “tight”, viz. , “तंग” and “कसा” were used only occasionally by the participants even after they were specifically told to adhere to Hindi in their narrations. This was the same with Hindi speaking population which had very less contact with English language.

| State Profiled: [loose] corpus | | | | | | | |
|--------------------------------|----|-------------------------------|-------------------------------|----------|--------|---|---------|
| Term | fr | $p\left(\frac{w}{C_T}\right)$ | $p\left(\frac{w}{C_L}\right)$ | f_{rL} | $p(w)$ | $\frac{p\left(\frac{w}{C_T}\right)}{p\left(\frac{w}{C_L}\right)}$ | |
| Without Stemming | | | | | | | |
| loose | 1 | 0.01190 | 4 | 0.00877 | 5 | 0.00209 | 4.61403 |
| ढर | 1 | 0.00190 | 4 | 0.00877 | 5 | 0.00509 | 4.61403 |
| ढरफ | 2 | 0.00380 | 7 | 0.01535 | 9 | 0.00916 | 4.03728 |
| compensation | 1 | 0.00190 | 3 | 0.00657 | 4 | 0.00407 | 3.46052 |
| peg | 1 | 0.00190 | 3 | 0.00657 | 4 | 0.00407 | 3.46052 |
| size | 1 | 0.00190 | 3 | 0.00657 | 4 | 0.00407 | 3.46052 |
| ढरफ | 2 | 0.00380 | 6 | 0.01317 | 8 | 0.00814 | 3.46052 |
| ढर | 1 | 0.00190 | 3 | 0.00657 | 4 | 0.00407 | 3.46052 |
| ढरफ | 2 | 0.00380 | 6 | 0.01317 | 8 | 0.00814 | 3.46052 |
| ढरफ | 3 | 0.00570 | 7 | 0.01923 | 10 | 0.01018 | 2.69152 |
| With Stemming | | | | | | | |
| loose | 1 | 0.01190 | 4 | 0.00877 | 5 | 0.00209 | 4.61403 |
| ढर | 1 | 0.00190 | 3 | 0.00657 | 4 | 0.00407 | 3.46052 |
| peg | 1 | 0.00190 | 3 | 0.00657 | 4 | 0.00407 | 3.46052 |
| ढर | 1 | 0.00190 | 3 | 0.00657 | 4 | 0.00407 | 3.46052 |
| ढरफ | 3 | 0.00570 | 9 | 0.01923 | 12 | 0.01221 | 3.46052 |
| ढरफ | 3 | 0.00570 | 7 | 0.01535 | 10 | 0.01018 | 2.69152 |
| ढर | 3 | 0.00570 | 7 | 0.01535 | 10 | 0.01018 | 2.69152 |
| ढर | 1 | 0.00190 | 2 | 0.00438 | 3 | 0.00305 | 2.30701 |
| ढर | 1 | 0.00190 | 2 | 0.00438 | 3 | 0.00305 | 2.30701 |
| ढर | 2 | 0.00380 | 4 | 0.00877 | 6 | 0.00610 | 2.30701 |

Table 1: Hindi State Profiled [loose] corpus: Top 10 words by conditional ratio

| State Profiled: [tight] corpus | | | | | | | |
|--------------------------------|----|-------------------------------|-------------------------------|----------|--------|---|---------|
| Term | fr | $p\left(\frac{w}{C_T}\right)$ | $p\left(\frac{w}{C_L}\right)$ | f_{rL} | $p(w)$ | $\frac{p\left(\frac{w}{C_T}\right)}{p\left(\frac{w}{C_L}\right)}$ | |
| Without Stemming | | | | | | | |
| tight | 9 | 0.01111 | 1 | 0.00219 | 10 | 0.01018 | 7.80228 |
| ढर | 7 | 0.01333 | 1 | 0.00219 | 8 | 0.00814 | 6.08844 |
| ढर | 5 | 0.00950 | 1 | 0.00219 | 6 | 0.00610 | 4.33160 |
| ढरफ | 4 | 0.00760 | 1 | 0.00219 | 5 | 0.00509 | 3.67688 |
| movement | 3 | 0.00570 | 1 | 0.00219 | 4 | 0.00407 | 2.60076 |
| ढरफ | 3 | 0.00570 | 1 | 0.00219 | 4 | 0.00407 | 2.60076 |
| ढर | 6 | 0.01140 | 2 | 0.00438 | 8 | 0.00814 | 2.60076 |
| ढर | 5 | 0.00950 | 2 | 0.00438 | 7 | 0.00712 | 2.16730 |
| ढर | 5 | 0.00950 | 2 | 0.00438 | 7 | 0.00712 | 2.16730 |
| ढरफ | 16 | 0.02041 | 7 | 0.01535 | 23 | 0.02321 | 1.98153 |
| With Stemming | | | | | | | |
| tight | 9 | 0.01111 | 1 | 0.00219 | 10 | 0.01018 | 7.80228 |
| ढरफ | 4 | 0.00760 | 1 | 0.00219 | 5 | 0.00509 | 3.46768 |
| movement | 3 | 0.00570 | 1 | 0.00219 | 4 | 0.00407 | 2.60076 |
| ढर | 3 | 0.00570 | 1 | 0.00219 | 4 | 0.00407 | 2.60076 |
| ढर | 6 | 0.01140 | 2 | 0.00438 | 8 | 0.00814 | 2.60076 |
| ढर | 5 | 0.00950 | 2 | 0.00438 | 7 | 0.00712 | 2.16730 |
| ढर | 5 | 0.00950 | 2 | 0.00438 | 7 | 0.00712 | 2.16730 |
| ढरफ | 18 | 0.03422 | 8 | 0.01535 | 26 | 0.02647 | 1.95967 |
| pele | 6 | 0.01140 | 3 | 0.00627 | 9 | 0.00916 | 1.73384 |

Table 2: Hindi Action Profiled [tight] corpus: Top 10 words by conditional ratio

| Action Profiled: [loose] corpus | | | | | | | |
|---------------------------------|----|-------------------------------|-------------------------------|----------|-------------|---|---------|
| Term | fr | $p\left(\frac{w}{C_T}\right)$ | $p\left(\frac{w}{C_L}\right)$ | f_{rL} | $p(w)$ | $\frac{p\left(\frac{w}{C_T}\right)}{p\left(\frac{w}{C_L}\right)}$ | |
| Without Stemming | | | | | | | |
| loose | 0 | 0.01780 | 10 | 0.00754 | Large Value | | |
| ढरफ | 1 | 0.00133 | 3 | 0.00519 | 4 | 0.00301 | 3.88235 |
| ढर | 1 | 0.00133 | 3 | 0.00519 | 4 | 0.00301 | 3.88235 |
| ढर | 1 | 0.00133 | 3 | 0.00519 | 4 | 0.00301 | 3.88235 |
| ढर | 3 | 0.00401 | 7 | 0.01211 | 10 | 0.00754 | 3.01960 |
| ढर | 4 | 0.00534 | 8 | 0.01584 | 12 | 0.00904 | 2.58823 |
| ढर | 1 | 0.00133 | 2 | 0.00346 | 3 | 0.00226 | 2.58823 |
| ढरफ | 1 | 0.00133 | 2 | 0.00346 | 3 | 0.00226 | 2.58823 |
| ढर | 1 | 0.00133 | 2 | 0.00346 | 3 | 0.00226 | 2.58823 |
| With Stemming | | | | | | | |
| loose | 0 | 0.01780 | 10 | 0.00754 | Very Large | | |
| ढर | 1 | 0.00133 | 6 | 0.01038 | 7 | 0.00602 | 7.76470 |
| ढर | 1 | 0.00133 | 5 | 0.00865 | 6 | 0.00431 | 6.47058 |
| ढर | 1 | 0.00133 | 5 | 0.00865 | 6 | 0.00431 | 6.47058 |
| ढर | 1 | 0.00133 | 4 | 0.00692 | 5 | 0.00359 | 5.17647 |
| ढर | 2 | 0.00267 | 6 | 0.01038 | 8 | 0.00574 | 3.88235 |
| ढर | 1 | 0.00133 | 3 | 0.00519 | 4 | 0.00287 | 3.88235 |
| ढरफ | 2 | 0.00267 | 5 | 0.00865 | 7 | 0.00602 | 3.23529 |
| ढर | 7 | 0.00865 | 17 | 0.02941 | 24 | 0.01724 | 3.12885 |
| ढर | 1 | 0.00133 | 2 | 0.00346 | 3 | 0.00215 | 2.58823 |

Table 3: Hindi Action Profiled [loose] corpus: Top 10 words by conditional ratio

| Action Profiled: [tight] corpus | | | | | | | |
|---------------------------------|----|-------------------------------|-------------------------------|----------|--------|---|---------|
| Term | fr | $p\left(\frac{w}{C_T}\right)$ | $p\left(\frac{w}{C_L}\right)$ | f_{rL} | $p(w)$ | $\frac{p\left(\frac{w}{C_T}\right)}{p\left(\frac{w}{C_L}\right)}$ | |
| Without Stemming | | | | | | | |
| tight | 10 | 0.01336 | 1 | 0.00173 | 11 | 0.00820 | 7.72727 |
| ढर | 7 | 0.00935 | 1 | 0.00173 | 8 | 0.00603 | 5.40600 |
| ढर | 5 | 0.00688 | 1 | 0.00173 | 6 | 0.00452 | 3.80361 |
| ढर | 9 | 0.01203 | 2 | 0.00346 | 11 | 0.00820 | 3.47727 |
| ढर | 4 | 0.00534 | 1 | 0.00173 | 5 | 0.00377 | 3.09090 |
| ढर | 3 | 0.00401 | 1 | 0.00173 | 4 | 0.00301 | 2.31818 |
| ढर | 3 | 0.00401 | 1 | 0.00173 | 4 | 0.00301 | 2.31818 |
| ढर | 3 | 0.00401 | 1 | 0.00173 | 4 | 0.00301 | 2.31818 |
| ढर | 3 | 0.00401 | 1 | 0.00173 | 4 | 0.00301 | 2.31818 |
| ढर | 5 | 0.00688 | 2 | 0.00346 | 7 | 0.00527 | 1.91811 |
| With Stemming | | | | | | | |
| tight | 10 | 0.01336 | 1 | 0.00173 | 11 | 0.00820 | 7.72727 |
| ढर | 3 | 0.00401 | 1 | 0.00173 | 4 | 0.00287 | 2.31818 |
| ढर | 3 | 0.00401 | 1 | 0.00173 | 4 | 0.00287 | 2.31818 |
| size | 7 | 0.00935 | 3 | 0.00519 | 10 | 0.00718 | 1.80360 |
| ढर | 11 | 0.01203 | 5 | 0.00865 | 16 | 0.01149 | 1.70909 |
| ढर | 6 | 0.00865 | 3 | 0.00519 | 9 | 0.00646 | 1.54545 |
| ढर | 2 | 0.00267 | 1 | 0.00173 | 3 | 0.00215 | 1.54545 |
| ढर | 2 | 0.00267 | 1 | 0.00173 | 3 | 0.00215 | 1.54545 |
| ढर | 2 | 0.00267 | 1 | 0.00173 | 3 | 0.00215 | 1.54545 |
| ढर | 2 | 0.00267 | 1 | 0.00173 | 3 | 0.00215 | 1.54545 |

Table 4: Hindi Action Profiled [tight] corpus: Top 10 words by conditional ratio

References

- [Madan Dabberu, Amitabha Mukherjee]. "Using Symbol Emergence to Discover Multi-Lingual Translations in Design". Proceedings of the ASME 2010 International Design Engineering Technical Conferences & Computers and Information in Engineering Conference Proceedings of IDETC/DTM 2010, 2010.
- [Dabberu, Madan Mohan and Mukerjee, Amitabha]. "Learning concepts and language for a baby designer". Design Computing and Cognition'10, 2011.
- [S V P Gopi Srinath, Nikhil Joshi, Prabhat Mudgal, Amitabha Mukerjee]. "Learning grounded semantics of Hindi nouns from video surveillance and user commentary". Proceedings of ICON-2010: 8th International Conference on Natural Language Processing, 2010.
- [Madan Dabberu, Amitabha Mukerjee]. "Computational models of tacit knowledge."
- [De Saussure, Ferdinand]. "Nature of the linguistic sign". Course In General Linguistics, 1916.