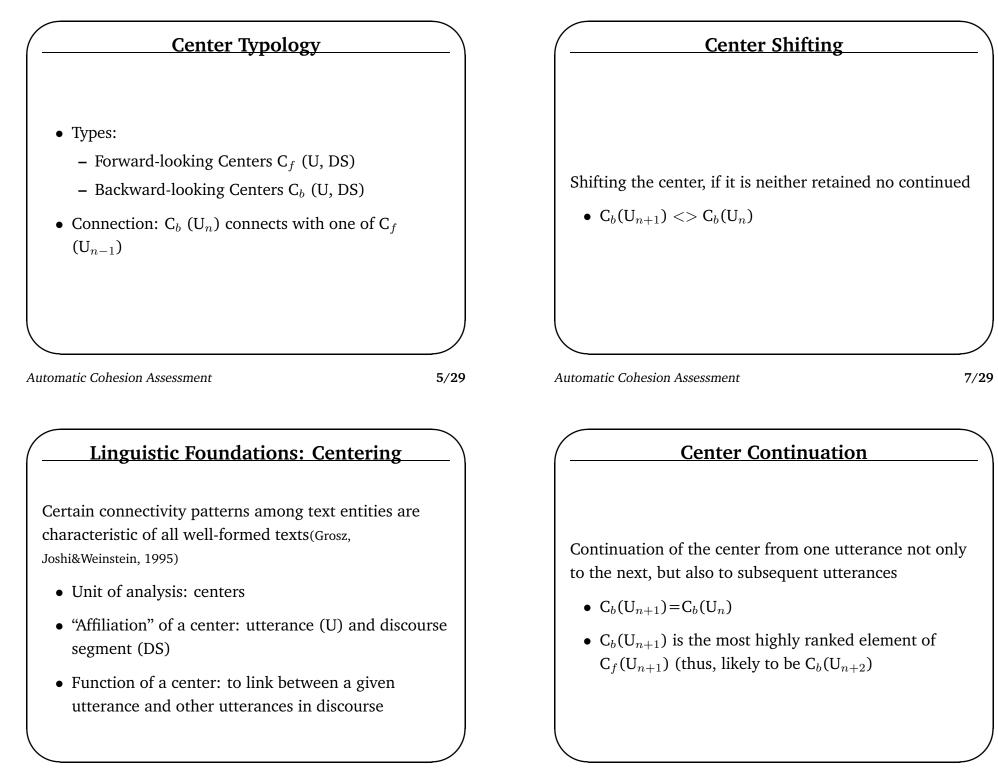
| Human vs  |   | Linguistic Foundations: Cohesion   |
|---|---|--|
|   |   | Cohesion: language devices that connect individual sentences into a  |
| A   | В   | unified whole  |
| In December 1988 a Pan Am jet was blown up over Lockerbie, Scotland,      | Secretary-General Kofi Annan said<br>Wednesday that he may travel to Libya  | Cohesion devices: repetition, coreference, ellipsis  |
| killing 270. Since 1992 Libya has   | next week in hopes of closing a deal  | 1. There was once a little girl and a little boy and a dog   |
| been under U.N. sanctions in effect                                       | to try two Libyan suspects in the Pan                                       | 2. And the sailor was their daddy  |
| intil the suspects are turned over<br>to United States or Britain. In Au- | Am Lockerbie bombing. The sanctions,<br>were imposed to force Libyan leader | 3. And the little doggy was white  |
| gust 1998 United States and Britain                                       | Moammar Gadhafi to turn the men   | 4. And they like the little doggy  |
| proposed a Netherlands trial. Libya<br>asked for guarantees that the sus- | over. Louis Farrakhan, the leader of a U.S. Muslim group, congratulated on  | 6. And they fed it   |
| pects would be incarcerated in  | his recovery from a hip injury.   | 7. And they ran away   |
| Libya. Kofi Annan planned a De-<br>cember 1988 Libyan trip to move        |   | 8. And then daddy had to go on a ship  |
| negotiations.   |   | 9. And the children misssed 'em  |
| matic Cohesion Assessment   | 1/29  | Automatic Cohesion Assessment  |
| matic Cohesion Assessment   | 1/29  | Automatic Cohesion Assessment  |
|   | 1/29  | Readability Models         • Goal: induce a model that can predict the degree text "well-formedness"   |
| Automatic Cohe  |   | • Goal: induce a model that can predict the degree   |
| Automatic Cohe  | esion Assessment  | Readability Models         • Goal: induce a model that can predict the degree text "well-formedness"         • Applications: summarization, question-answering                                       |
| Regina  | esion Assessment  | <ul> <li>Readability Models</li> <li>Goal: induce a model that can predict the degree text "well-formedness"</li> <li>Applications: summarization, question-answering machine-translation</li> </ul> |



### Discussion on Centering

- Until now: always based on manual annotations
- Never used in applications

Does it really work?

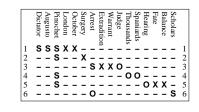
Automatic Cohesion Assessment

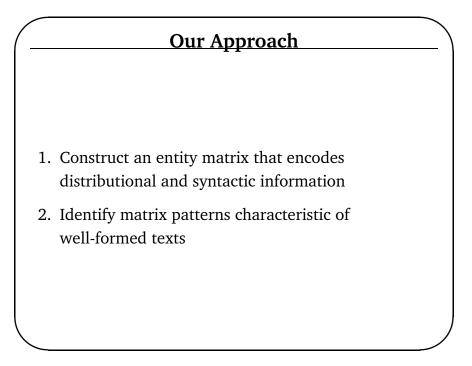
9/29

| / | Coherent   | Discourse   |
|---|--|---|
|   |  |   |
|   | Coherence is established via                           | center continuation                                     |
|   | John went to his favorite music store to buy a piano.  | John went to his favorite music store to buy a piano.   |
|   | He had frequented the store for many years.            | It was a store John had fre-<br>quented for many years. |
|   | He was excited that he could fi-<br>nally buy a piano. | He was excited that he could fi-<br>nally buy a piano.  |
|   | He arrived just as the store was closing for the day.  | It was closing just as John ar-<br>rived.               |
|   |  |   |

### Entity matrix

- 1. [Former Chilean dictator Augusto Pinochet]<sub>S</sub>, was arrested in  $[London]_X$  on  $[14 \text{ October}]_X$  1998.
- 2. [Pinochet]<sub>S</sub>, 82, was recovering from [surgery]<sub>X</sub>.
- 3. [The arrest]<sub>S</sub> was in [response]<sub>x</sub> to [an extradition warrant]<sub>x</sub> served by [a Spanish judge]<sub>O</sub>.
- 4. [Pinochet]  $_{\rm S}$  was charged with murdering [thousands]\_0, including many [Spaniards]\_0.
- 5. [Pinochet]<sub>S</sub> is awaiting [a hearing]<sub>O</sub>, [his fate]<sub>X</sub> in [the balance]<sub>X</sub>.
- 6. [American scholars]  $_{s}$  applauded the [arrest] $_{o}$ .





| istribution of syntactic tags |
|-------------------------------|
| istribution of syntactic tags |
| HRS MRS LRS                   |
| s s 0.020 0.014 0.010         |
| s o 0.012 0.005 0.004         |
| 0.417 0.433 0.450             |
|                               |
| Computation of Entity matrix  |
|                               |
|                               |
|                               |

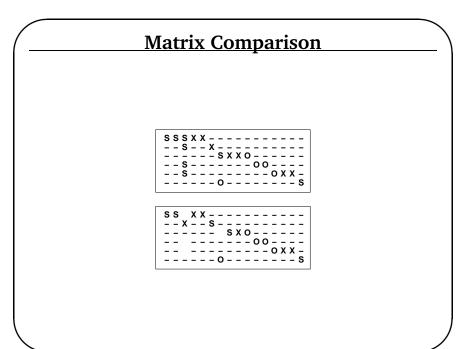
### **Transformations**

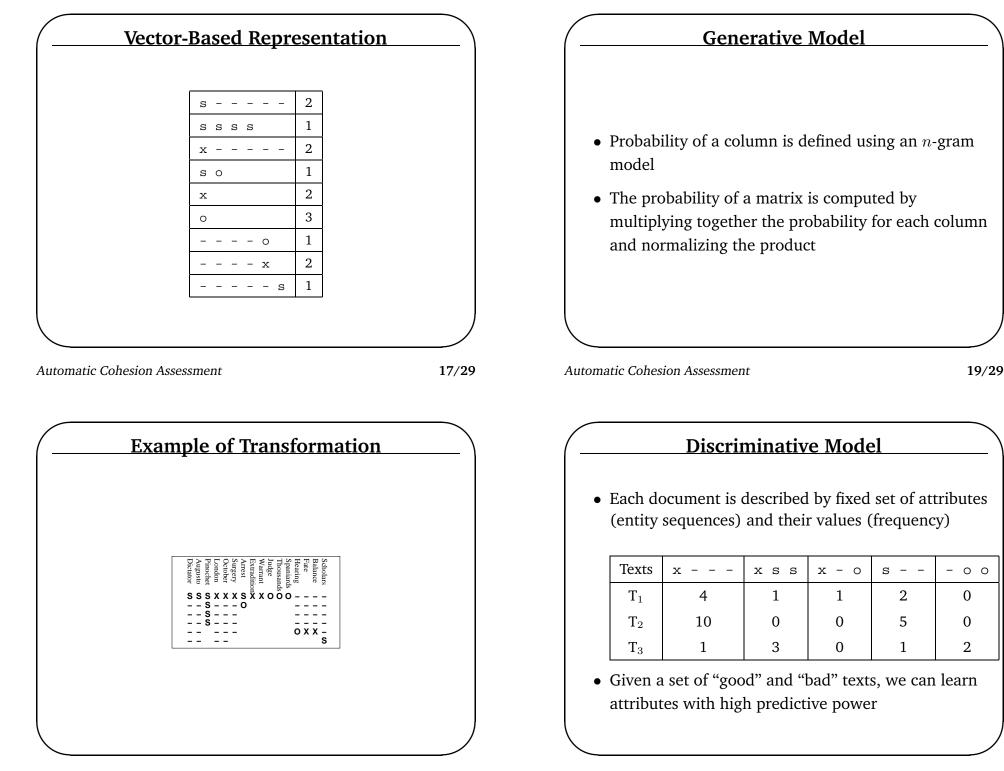
Goal: reduce the variability in matrix representation

| L  | Original | Transformed |
|----|----------|-------------|
| 1. | - s o    | s o         |
| 2. | - s o    | s o         |
| 3. | - s o    | S O         |

Automatic Cohesion Assessment

15/29





### Experiments: Data

| Humans | HRS  | MRS  | LRS  |
|--------|------|------|------|
| 5.13   | 4.42 | 4.32 | 3.60 |

#### Results of Anova Analysis:

- Human summaries are more cohesive than machine generated ones
- HRS is not significantly different from MRS
- Both HRS and MRS are significantly more cohesive than LRS

Automatic Cohesion Assessment

21/29

### **Experiments: Data**

- Data: Outputs of three multi-document summarization systems that participated in DUC'2003 and corresponding human summaries
  - High grammaticality scores
  - Variability in readability scores: High (HRS), Medium (MRS) and Low (LRS)
  - Overall 64 summaries
- Procedure: the judge assigns readability score on a seven point scale
  - 183 summaries (23 people per summary)

### **Model Comparison: Baselines**

- Readability Measures a function of the average sentence length and the average number of syllables (Flesh, 1951)
- Word-based Models the average word overlap of adjacent sentences (Foltz&Kintsch&Landauer, 1998)
- Vector-based Models the average distances between adjacent sentences based on word distributional properties (Foltz&Kintsch&Landauer, 1998)
- Taxonomy-based Models the average distances between adjacent sentences based on WordNet (Lin, Resnik)

Automatic Cohesion Assessment

23/29

# Agreement

Why not to use kappa?

- Function: Upper-bound on human performance
- Procedure: Leave-one-out resampling (Weiss&Kulikowski)
- Result: Agreement = .612 (Min = .107, Max = .975, SD = .230)

### **Results: Generative Model**

Correlation between human rating and the models

| Model                    | Correlation |
|--------------------------|-------------|
| Flesh Readability Index  | .010        |
| Word-based Model         | .113        |
| Latent Semantic Analysis | .184        |
| Taxonomy-based (Lin)     | 125         |
| Taxonomy-based (Resnik)  | 176         |
| Entity Matrix            | .314**      |
| *p < .05 (2-taile        | d)          |
| **p < .01 (2-taile       | ed)         |

Automatic Cohesion Assessment

25/29

### Generative Model: Implementation

- Applied to 6-letter alphabet at various level of compression
- Trained on DUC human summaries
- Tested on machine summaries

### **Results: Discriminative Model**

| Trans | Vec. Size | 2-way | 3-way |
|-------|-----------|-------|-------|
| Base  | -         | 69%   | 37.5% |
| 0     | 354       | 73%   | 53.3% |
| 1     | 101       | 97%   | 59.4% |
| 2     | 88        | 97%   | 59.4% |
| 3     | 77        | 97%   | 64.1% |
| 4     | 73        | 78%   | 64.1% |
| N     | 56        | 73%   | 62.5% |

Automatic Cohesion Assessment

27/29

## Discussion: Generative Model

- No correlation for traditional cohesion model due to redundancy
- High negative correlation for Wordnet-based models!
- Best results on the tranformation 3

### **Future directions**

- Dependence on genre
- Contribution of different linguistic features
  - Preliminary results: anaphora doesn't help
- More sophisticated model (unsupervised grammar induction, gap modeling)

Automatic Cohesion Assessment

29/29

### Discussion: Discriminative Model

- Most predictive patterns: [s x], [x o], [s s] and [s s s]
- Baselines: binary 67%, trinary 37,5%
- Transformation 3 is optimal in all the cases