# Predicting MT Fluency from
# IE Precision and Recall

Tony Hartley, Brighton, UK

Martin Rajman, EPFL, CH

# ISLE Keywords

- Quality of output text as a whole
  - Fluency (itself as a predictor of fidelity)
  - Utility of output for IE
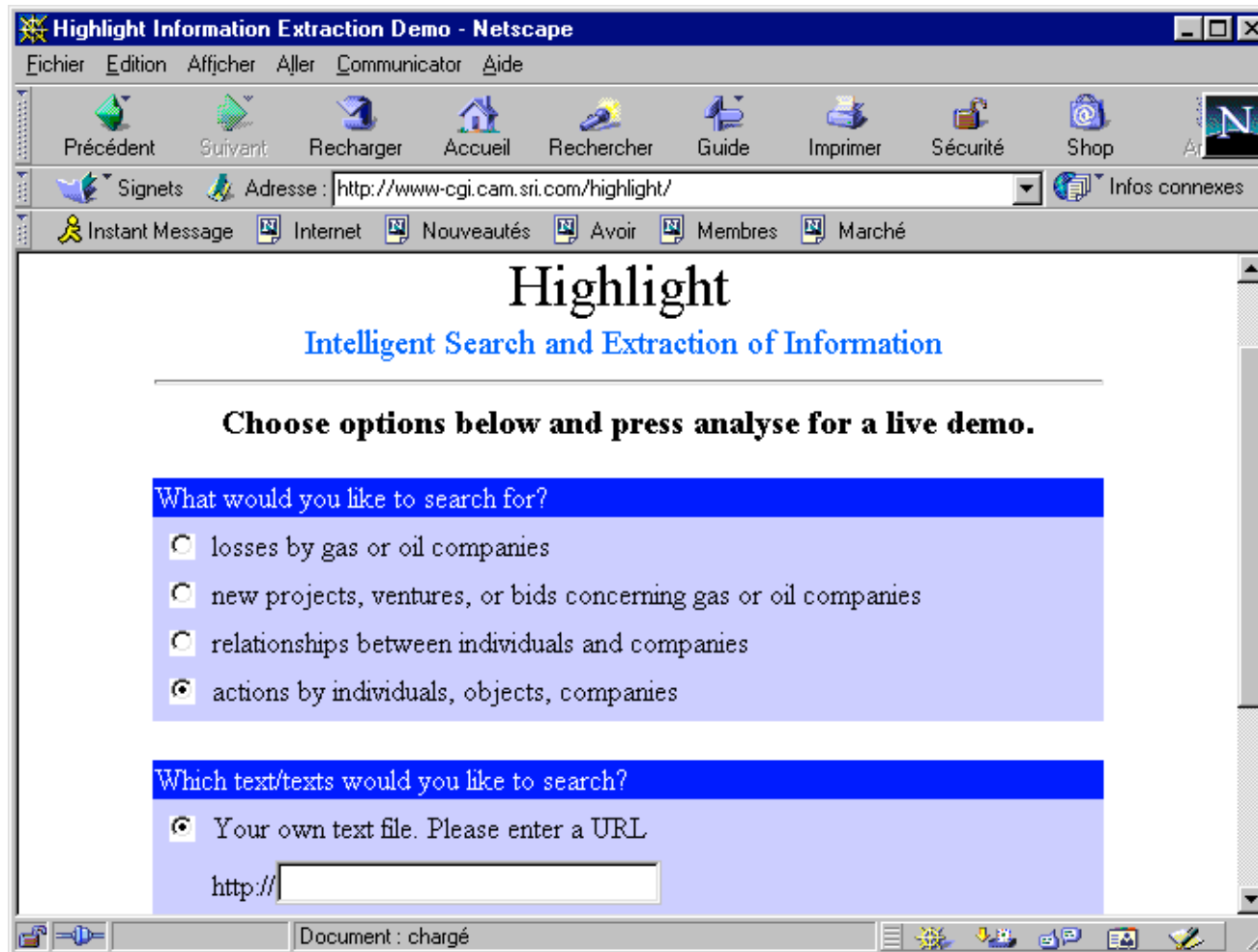- Automation
- F-score, Precision, Recall,

# Data Sources

- Fluency
  - DARPA 94 scores for French => English
  - 100 source texts * (1 HT + 5 MT)
- IE
  - Output from SRI *Highlight* engine
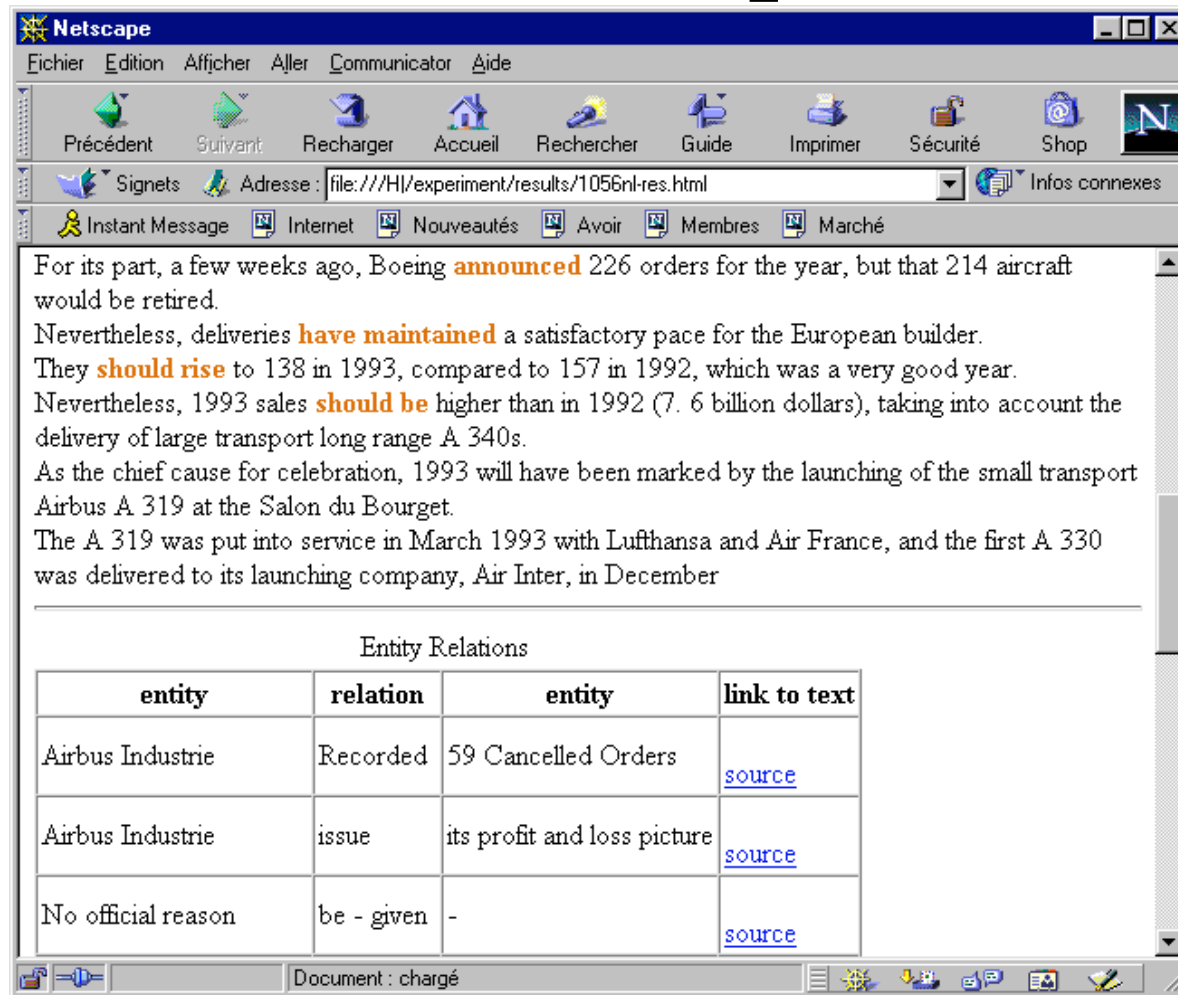  - `http://www-cgi.cam.sri.com/highlight/`

# Data Sampling

- Select 3 'systems'
  - Human expert -- mean fluency: 0.851934
  - S1 -- highest mean fluency: 0.508431
  - S2 -- lowest mean fluency: 0.124882
- Cover 'best to worse' for each
  - Aim -- select 20 texts
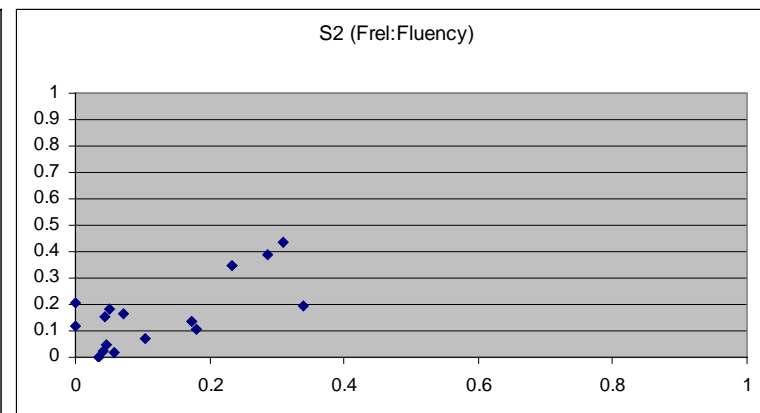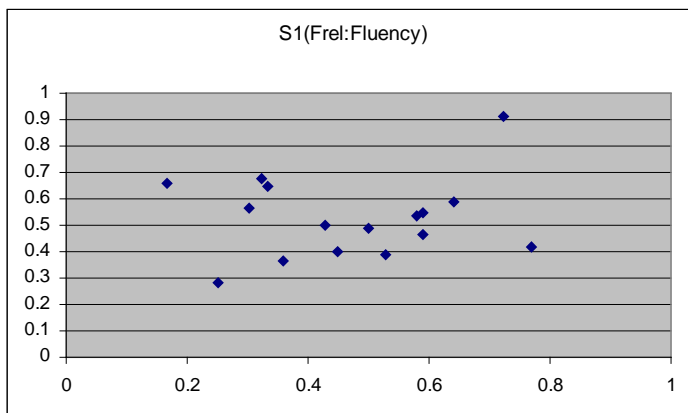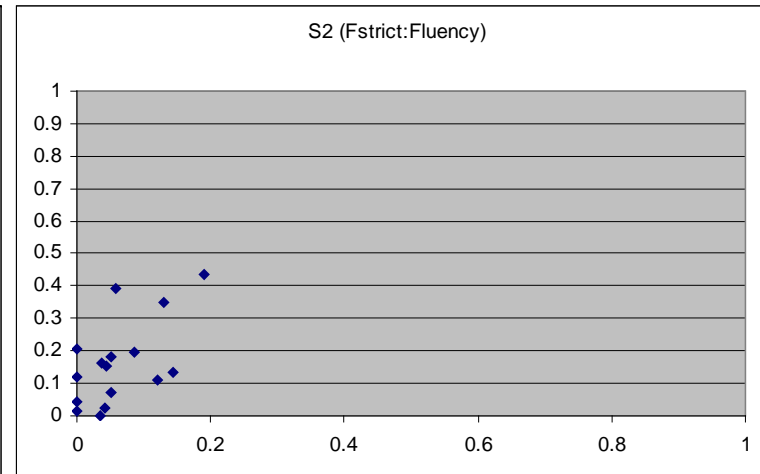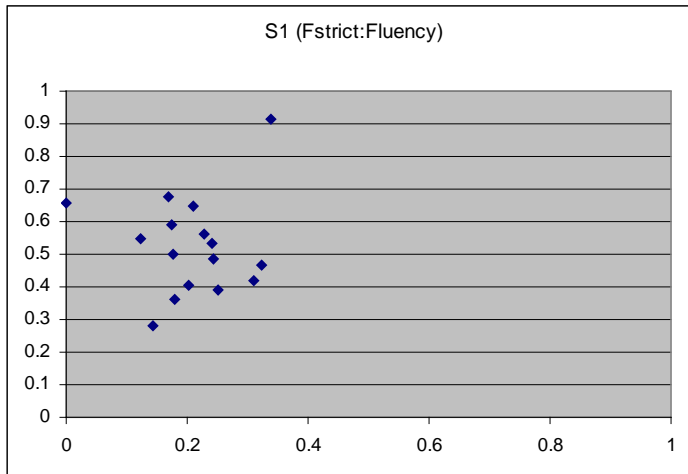  - Actual -- 17 texts

# IE Task

# IE Output

# IE Scoring

- Manual, but using simple, objective rules
- 'Gold Standard' given by expert translation
  - number of cells filled
- 'Strict' match = identity
- 'Relaxed' match = near miss
  - *quarter :: third school quarter*
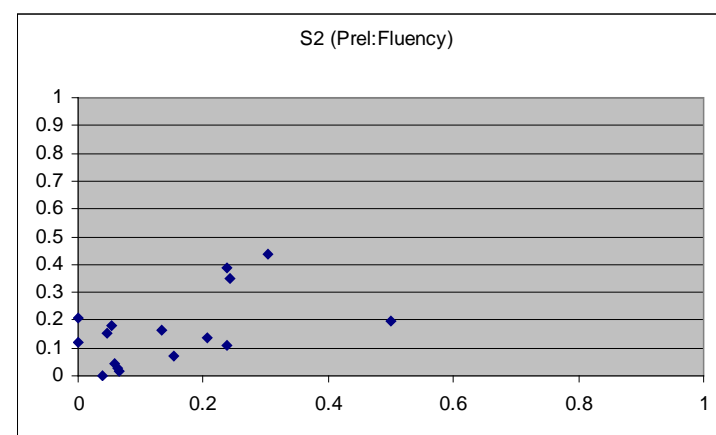  - *hoped :: wished*
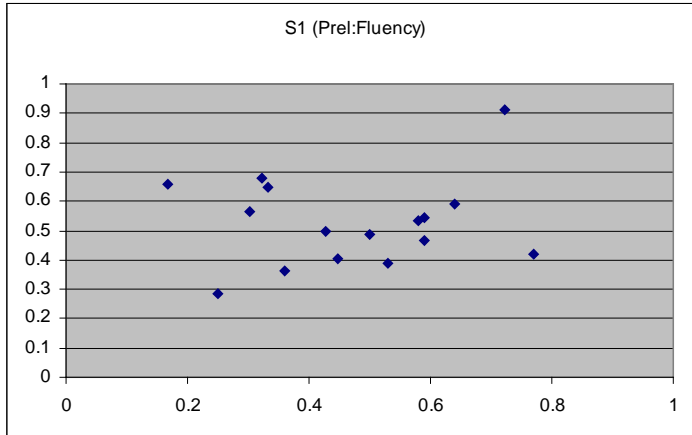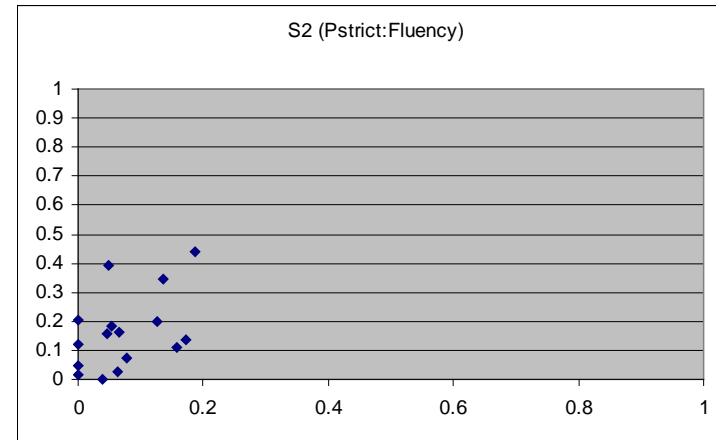  - *Airbus Industrie :: Airbus industry*

# Precision and Recall

- Precision *P*
  - (correct cells / actual cells)

- Recall *R*
  - (actual cells / reference cells)
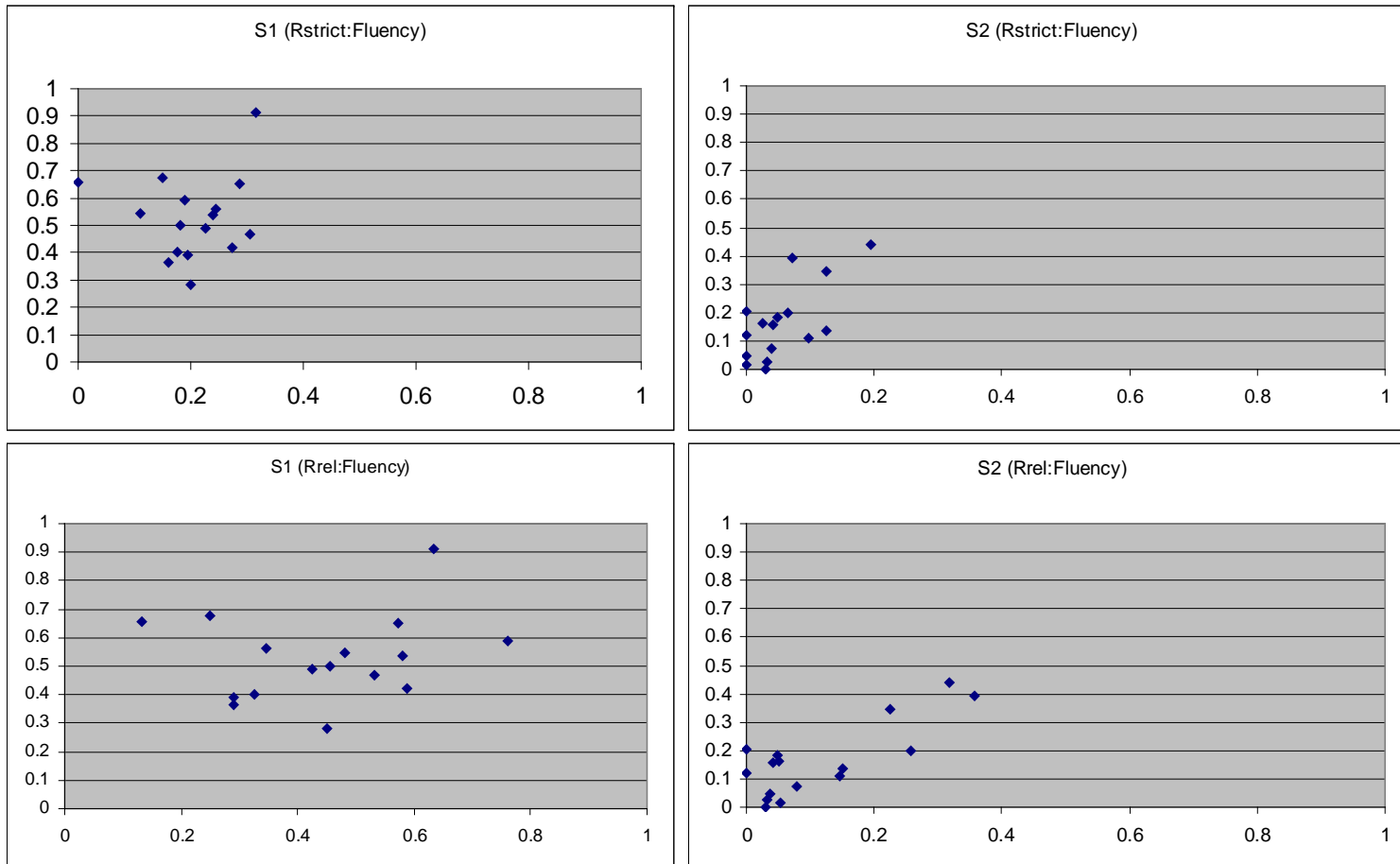
- F-score
  - 2 * P * R / (P + R)

# Correlations (not): F-score

# Correlations (not): Precision

# Correlations (not): Recall

# Conclusions

- A significant correlation is observed for S2 (for precision and recall)

- The same is not true for S1!…

- For S1, the correlation is higher for the F score than for P and R scores

# Observations on translation data

- MT yields structure preservation by default
- Does 'free-ness' of expert translation tarnish Gold Standard?
  - Finite verb in ST => non-finite verb in TT (IE under-retrieves)
  - Nominalisation in ST => finite verb in TT (IE over-retrieves)

# Observations on IE data

- Poor performance of *Highlight*
  - particularly on expert translations?
- Long distances between Entities and Relator
  - post-modification of NP
  - parenthetical material

# Next Steps

- Correlate with Adequacy rather than Fluency?

- Use whole data set (S2 bad choice?)

- Focus on extracted doubles and triples rather than single cells?

- Use another IE engine?

- Use ST to establish Gold Standard?