

NIST-ARPA Interagency Agreement: Human Language Technology Program

David S. Pallett, Principal Investigator

National Institute of Standards and Technology
Room A216, Building 225 (Technology)
Gaithersburg, MD 20899

PROJECT GOALS

1. To coordinate the design, development and distribution of speech and natural language corpora for the ARPA Spoken Language research community, and the use of these corpora for technology development and evaluation.
2. To design, coordinate the implementation of, and analyze the results of performance assessment benchmark tests for ARPA's speech recognition and spoken language understanding systems.
2. Participate in the ATIS SemEval effort, probably including the development of detailed test and reporting protocols for a "dry run" of an ATIS SemEval test.
3. Collect additional ATIS data at NIST as appropriate.
4. Continue to participate in the development of improved speech transcription and scoring procedures, in ATIS principles of Interpretation documents, and in cooperation with the annotators at SRI, in "bug-report adjudication".

RECENT RESULTS

1. Participated, with SRI International, in annotation and "bug fixes" for the ATIS MADCOW-collected corpora.
2. Installed BBN-developed and SRI-developed ATIS technology at NIST, and used this data to collect test and training data using subjects recruited from the Gaithersburg, MD area.
3. Produced speech corpora on recordable and pressed CD-ROM media in collaboration with the Linguistic Data Consortium.
4. Participated in discussions regarding implementation of the Semantic Evaluation (SemEval) glass box test protocols.
5. Prepared for, and implemented benchmark tests for the Wall Street Journal-based Continuous Speech Recognition (WSJ-CSR) corpus using the Hub-and-Spoke test paradigm and for the 46-city ATIS corpus.
5. Review the use of phonologically-motivated string alignment software for use in scoring speech recognition system output.
6. Prepare for and implement benchmark tests in the WSJ-CSR and ATIS domains in the November 1994 time frame.
7. Participate in the endeavors of the CCCC and MADCOW communities...

PLANS

1. Continue to collaborate with the LDC, its data collection and annotation contractors, and the MADCOW community with regard to data collection, annotation, screening and quality control procedures, and (as appropriate), to produce CD-ROMs for early release within the community of test participants.