# End-to-End Learning of Co-Speech Gesture Generation for Humanoid Robots
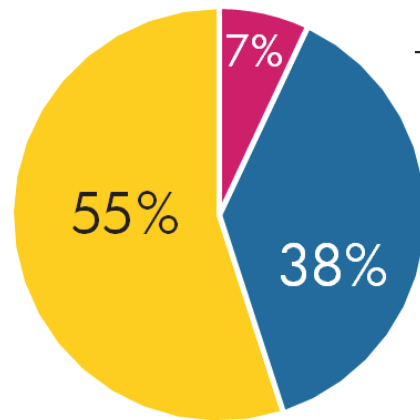
**Youngwoo Yoon**, Woo-Ri Ko, Minsu Jang, Jaeyeon Lee, Jaehong Kim

**ETRI** Electronics and Telecommunications Research Institute

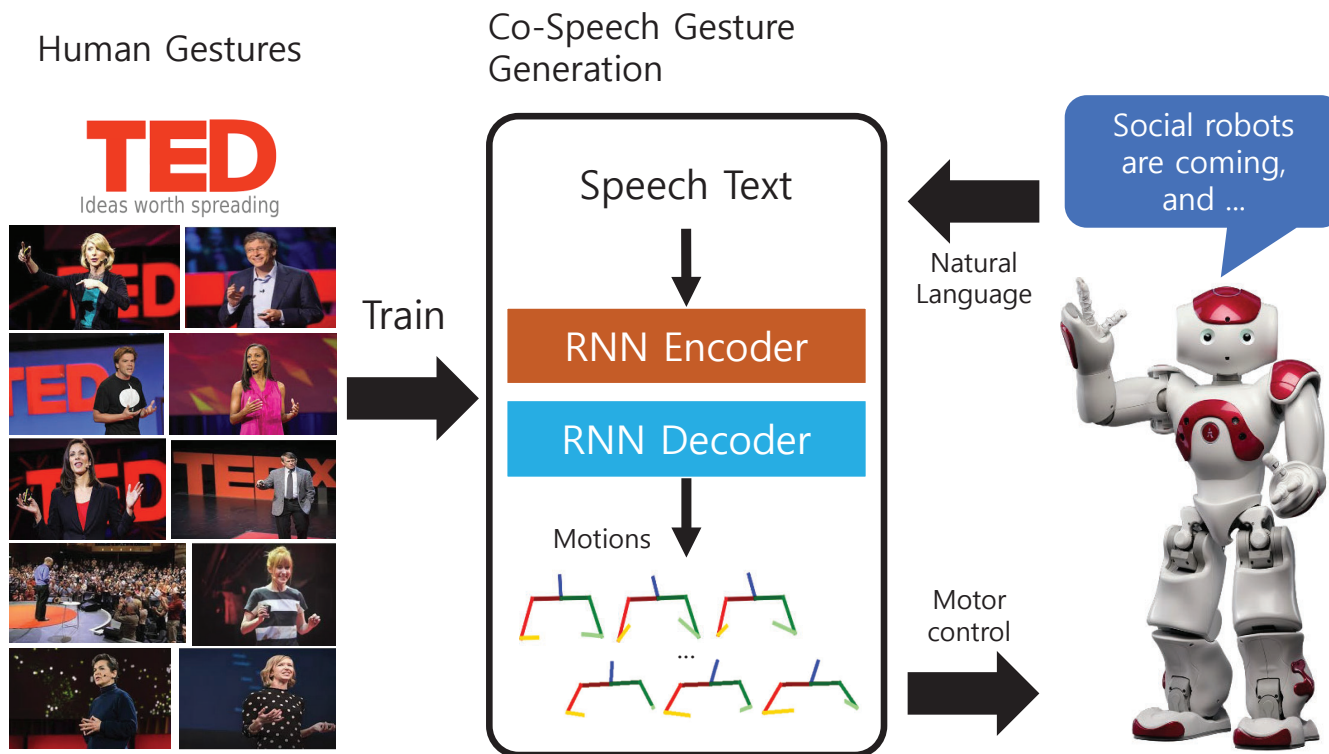# Motivation



Dr. Albert Mehrabian's 7-38-55% Rule

**Elements of Personal Communication**
- 7% spoken words
- 38% voice, tone
- 55% body language
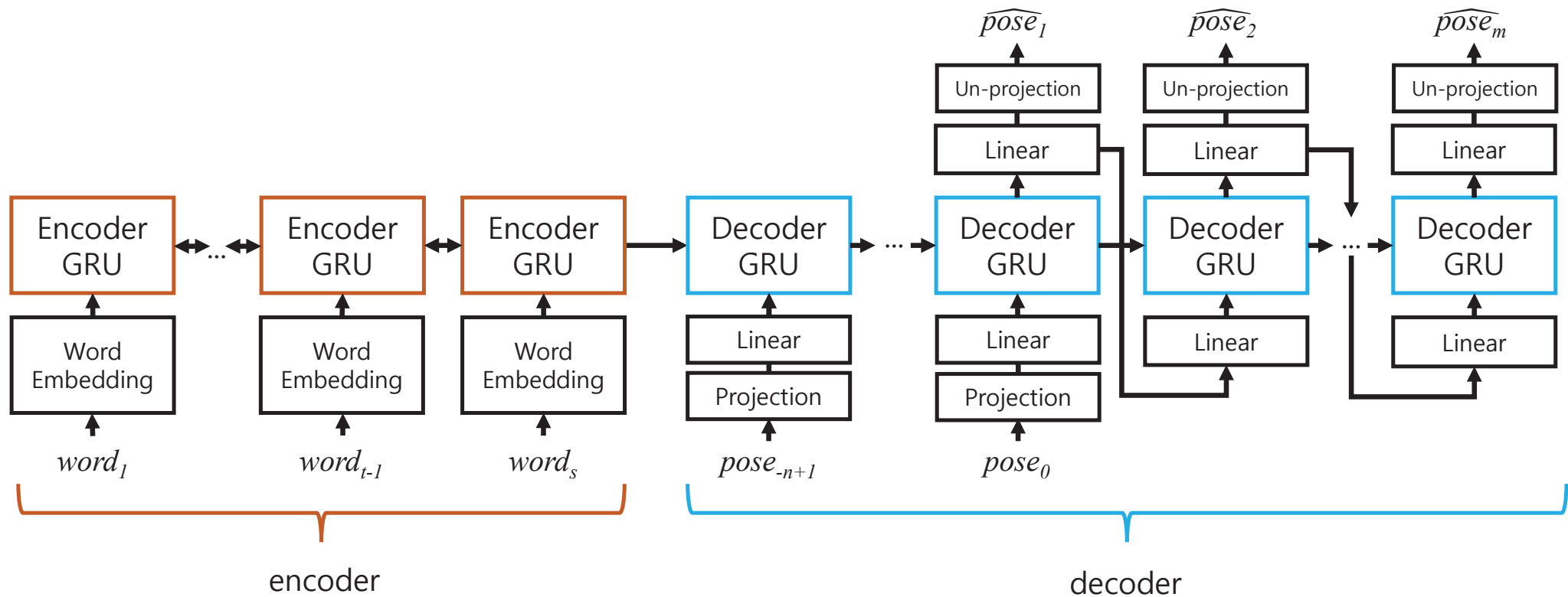
- Co-speech gesture
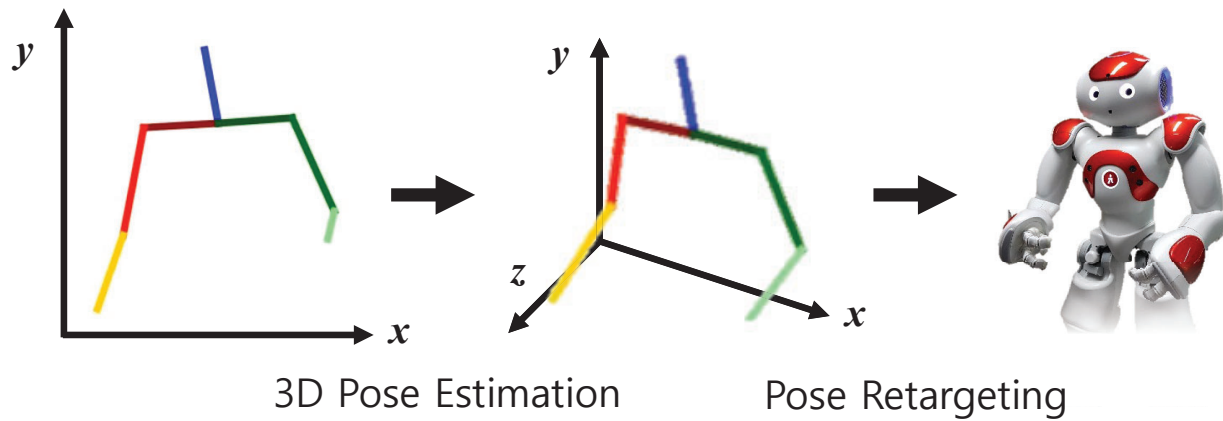  - Generating upper-body motion from speech text

# Overview



Human Gestures

Co-Speech Gesture Generation

Train

Speech Text

RNN Encoder

RNN Decoder

Motions

...

Natural Language

Motor control

Social robots are coming, and ...

# End-to-end Architecture

# Robot Prototype



3D Pose Estimation       Pose Retargeting

# Demo Video

# Thanks you

TED Dataset is available on

https://sites.google.com/view/youngwoo-yoon/projects/co-speech-gesture-generation