# Factor structure of a spoken Chinese test: investigating five subskill scores for diagnosis

LTRC 2013 in Seoul

Masanori Suzuki

Yujie Jia

Pearson Knowledge Technologies

# Overview of the Presentation

- Background on Spoken Chinese Test

- Test Structure

- Scoring logic

- Research Questions

  - Construct validation

  - Reliability analysis of subscores

- Analysis and Results

- Summary

# Spoken Chinese Test

- Fully automated test of spoken Mandarin Chinese
- Jointly developed between Peking Univ and Pearson
- Development period: June 2010 - August 2012
- Integrated listening-speaking item types
- Delivery options: phone, computer
- Scored automatically by in-house, custom-developed speech processing technologies and computerized scoring systems

# Test Design

# Spoken Chinese Test

- 8 item types
- 70 items
- 20 min
- Overall score
- 5 analytic subscores
- Score Scale: 20-80

# Test Construct

*Facility in spoken Chinese*

*The ability to understand spoken Chinese on everyday topics and to respond intelligibly in Chinese at a native-like conversational pace*



**Listen**

hear utterance
extract words
get phrase structure
decode propositions
contextualize
infer demand (if any)

articulate response
build clause structure
select lexical items
construct phrases
select register
decide on response

**Speak**

*Adapted from Levelt, 1989*

PEARSON

# Test Structure

| Part | Item Type | Response Characteristics | Scored Trait(s) |
|------|-----------|--------------------------|-----------------|
| A | Tone Phrases | Word, Phrase | Tone Production |
| B | Read Aloud | Sentence | Pronunciation, Tone production |
| C | Sentence Repetition | Sentence | Grammar, Pronunciation, Fluency, Tone production |
| D | Short Answer Questions | Word, Phrase | Vocabulary |
| E | Recognize Tones - Word | Word | Tone reception |
| F | Recognize Tone - Sentence | Word | Tone reception |
| G | Sentence Builds | Sentence | Grammar, Pronunciation, Fluency, Tone production |
| H | Passage Retellings | Passage | Vocabulary, Fluency |

# 2 Test Methods

- Test Method 1: Short Responses

    Tone Phrase

    Read Aloud

    Repeat

    Short Answer Questions

    Sentence Builds

    Recognize Tones


- Test Method 2: Long Responses

    Passage Retelling

# Example: Short response

Sentence Repeats

要下雨了。
*It's going to rain.*


后来他又去了一次。
*Afterwards, he went back again.*


报纸上说明天下午两点开始。
*The newspaper says that it starts at two o'clock tomorrow afternoon.*

*Scored traits: Grammar, Pron, Fluency, Tone Production*

# Example: Long Response

Passage Retelling

手机太好用了。不管你在哪儿，都能随时打电话，发短信。现在的手机功能更多，不但可以听音乐，还能上网呢。

*Cellphones are great.  No matter where you are, you can always make calls and send text messages.  Nowadays cellphones have even more functions.  You can use them not only to listen to music but also to surf the Internet.*

*Scored traits: Vocabulary, Fluency*

# Score reporting design



Five analytic subscores

# Multi-Trait, Multi-Method

Grammar | Vocabulary | Fluency | Pronunciation | Tone

Tone Phrase | Read Aloud | Repeat | Short Answer Question | R Tone Word | R Tone Sent | Sentence Build | Passage Retell

**20 minutes**

ALWAYS LEARNING

PEARSON

# Five Analytic Subscores

| Grammar | Vocabulary | Fluency | Pronunciation | Tone |
|---------|------------|---------|---------------|------|

- Fox and Fraser (2009)

  "What is surprising here is that the test developers have not emphasized the obvious diagnostic properties of the *The Versant Spanish™ Test.* " *(p.319)*

- SCT's test design and scoring logic is similar to that of Versant Spanish

  → *Can SCT's subscores be useful for diagnostic purposes?*

# Conditions for Being Useful

- SCT should exhibit evidence that the test measures facility in spoken Chinese (construct validation)

- SCT's subscores should be stable estimates (reliability)

# Research Questions

RQ 1: Does the SCT discriminate as expected according to known populations?

RQ 2: What is the factor structure of the SCT test scores?

RQ 3: Are SCT's subscores reliable enough to be taken as stable indicators of test-takers' strengths and weaknesses?

# Native Data Set

- 1,822 completed tests

- Various dialect groups

- Used to develop scoring systems

| Dialect | % |
|---------|-----|
| Mandarin | 57% |
| Wu | 13% |
| Min | 9% |
| Yue | 6% |
| Xiang | 3% |
| Gan | 2% |
| Hakka | 2% |
| Jing | 0.4% |

# Learner Data Set 1

- 3,845 completed tests
- Various countries and L1 backgrounds
- Used to develop scoring systems

| L1 | % |
|----|----|
| English | 17% |
| Japanese | 13% |
| Korean | 9% |
| Russian | 7% |
| Spanish | 6% |
| Arabic | 5% |
| Thai | 5% |
| Bengali | 5% |

# Learner Data Set 2

- 166 learners of Chinese

- Various countries and L1 backgrounds

- Set aside for validation

| L1 | % |
|----|----|
| English | 24% |
| Korean | 13% |
| Cantonese | 11% |
| Japanese | 8% |

# Research Questions

RQ 1: Does the SCT discriminate according to known populations as expected?

RQ 2: What is the factor structure of the SCT test scores?

RQ 3: Are SCT's subscores reliable enough to be taken as stable indicators of test-takers' strengths and weaknesses?

# Educated native Chinese speakers of different dialects

# *Educated native speakers and learners of Chinese*

# Educated native speakers, learners, heritage-Mandarin speakers, and heritage-NonMandarin speakers

# *Accuracy of SCT's Automated Scoring*

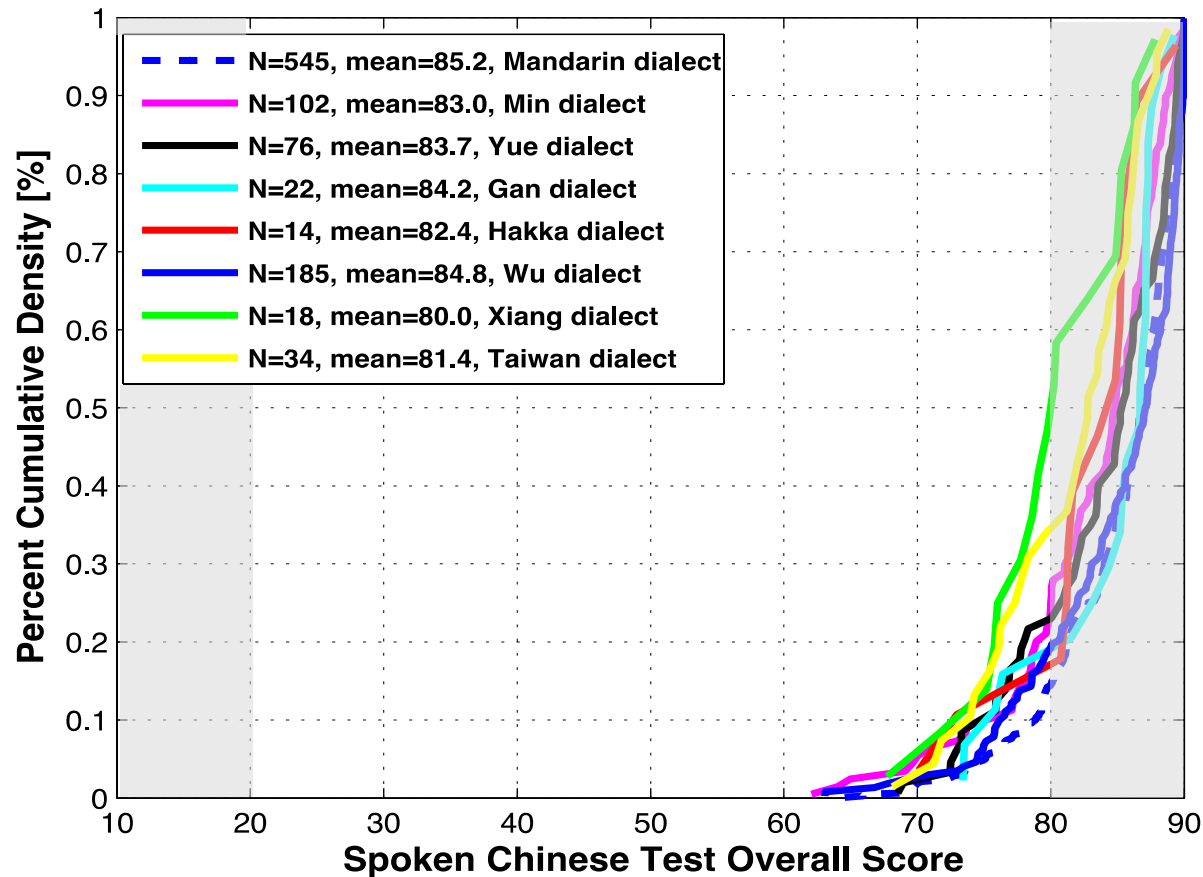| Score | Machine – Human Correlation (n=166) |
|---|---|
| Overall | **0.98** |
| Grammar | 0.97 |
| Vocabulary | 0.97 |
| Fluency | 0.93 |
| Pronunciation | 0.90 |
| Tone | 0.92 |

# Automated Scoring vs. Human Scoring

# Research Questions

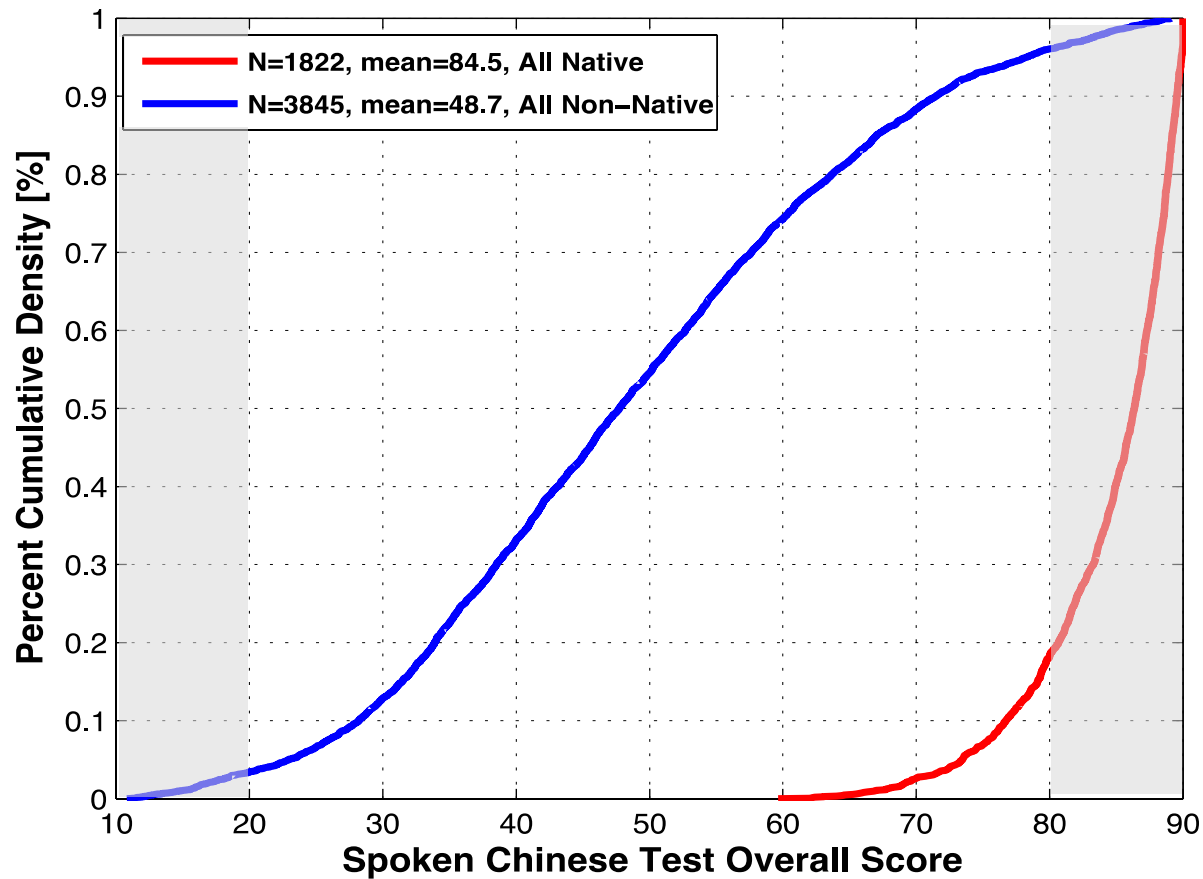RQ 1: Does the SCT discriminate according to known populations as expected?

RQ 2: What is the factor structure of the SCT test scores?

RQ 3: Are SCT's subscores reliable enough to be taken as stable indicators of test-takers' strengths and weaknesses?
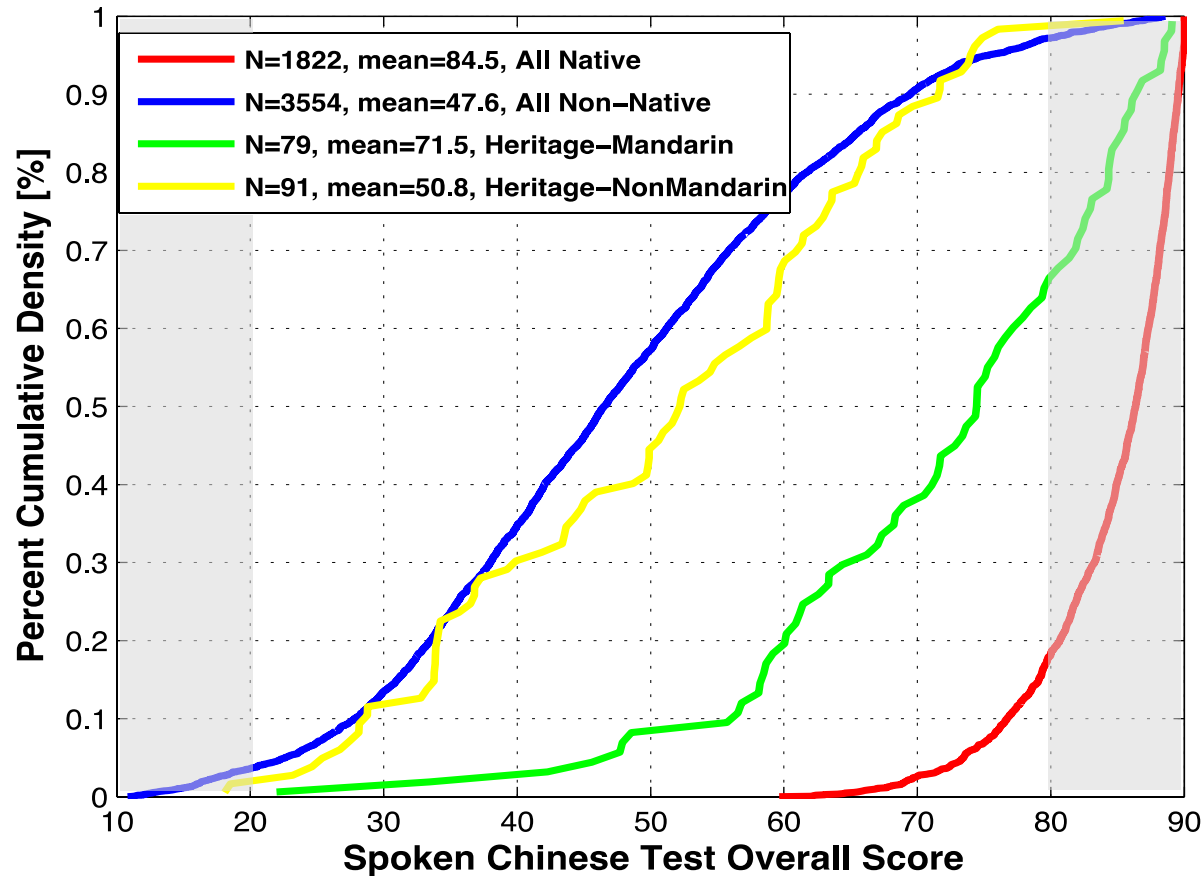
# Eight Variables

*1*

| Grammar | Vocabulary | Fluency | *6* Pronunciation | Tone |
|---------|-----------|---------|--------------------|------|

*2*  *3*  *4*  *5*  *7*  *8*

| ShortVoc | LongVoc | ShortFlu | LongFlu | Tone Production | Tone Recognition |
|----------|---------|----------|---------|-----------------|------------------|

# Data Analysis

- Confirmatory Factor Analysis

  Four CFA models were hypothesized and tested in this study

  1) One-factor model (Unidimensional construct)

  2) Correlated two-factor model

  3) Bi-factor model

  4) Correlated trait-uncorrelated method model

# One-Factor Model

# Correlated Two-Factor Model

# Bi-Factor Model

PEARSON

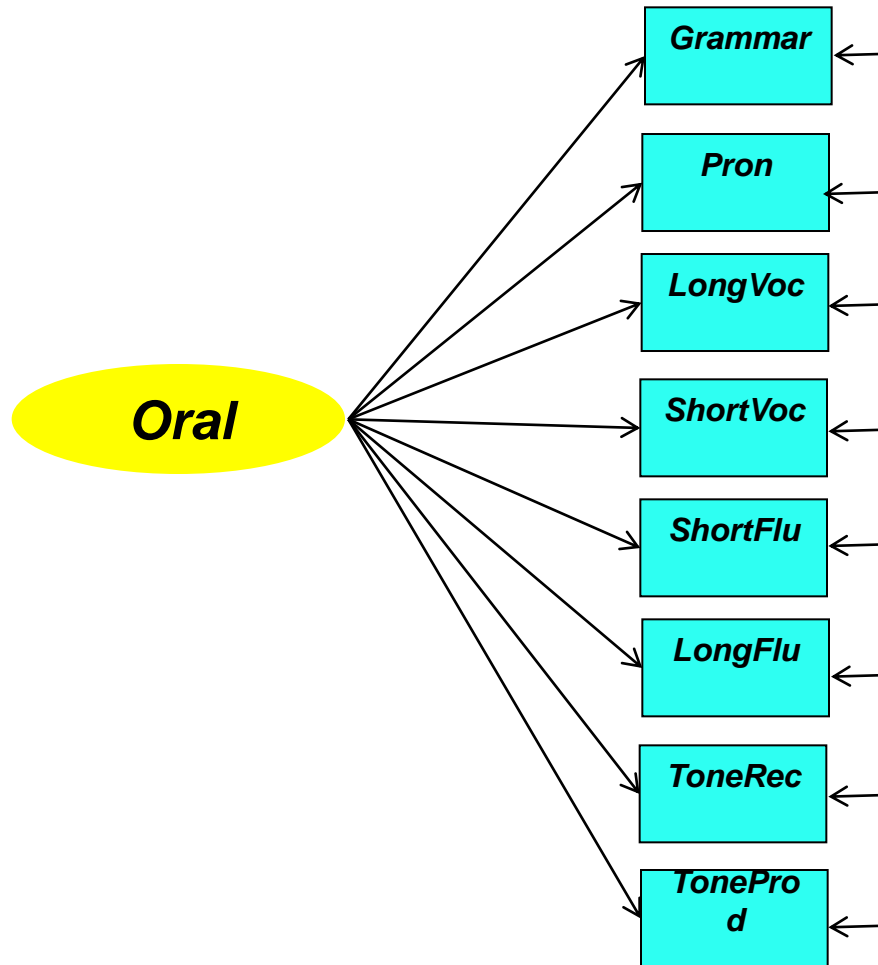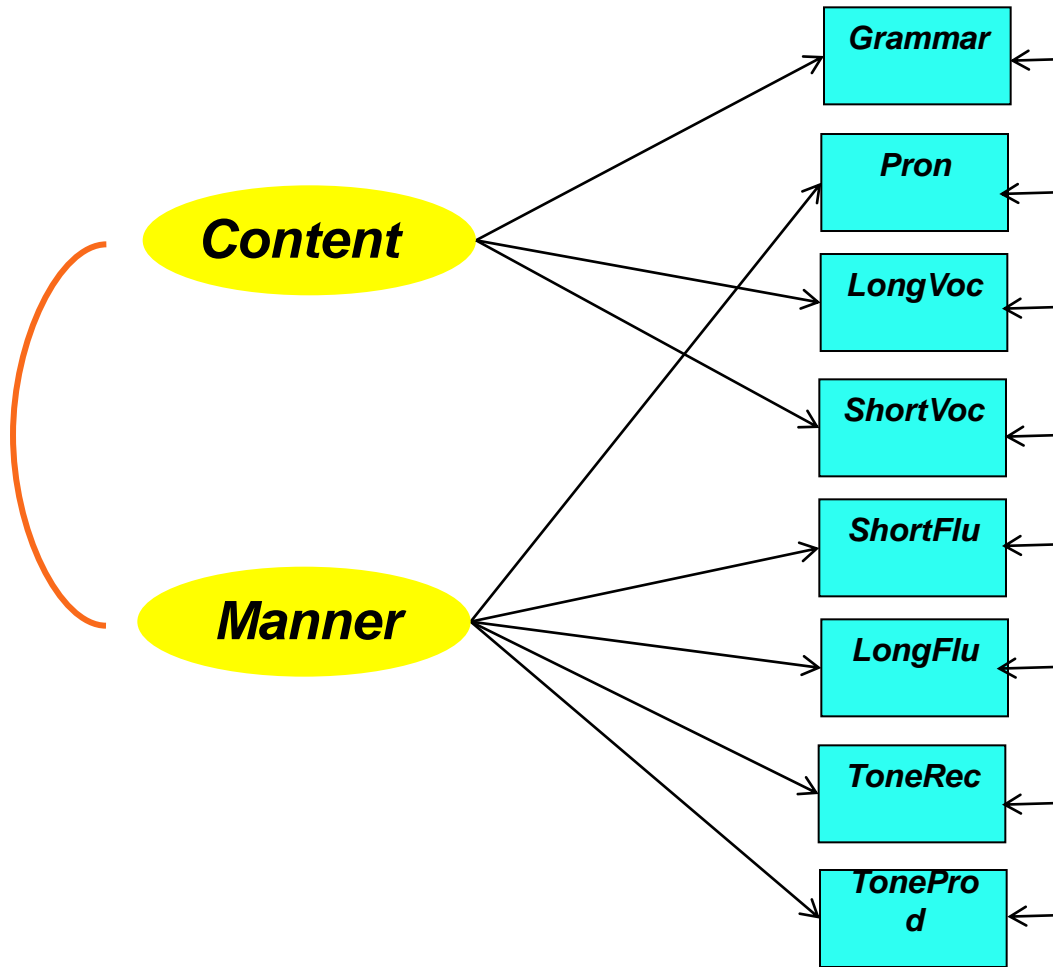# Correlated-Traits, Uncorrelated-Method

# Results

| Model Description | One-factor | Correlated two-factor | Bi-factor | Corrrelated trait-uncorrelated method |
|---|---|---|---|---|
| df | 20 | 19 | 16 | 15 |
| Minimum fit function chi-square | 3256.17 | 2151.25 | 336.62 | 194.35 |
| P value | <0.001 | <0.001 | <0.001 | <0.001 |
| RMSEA | 0.273 | 0.227 | 0.148 | 0.074 |
| CFI | 0.836 | 0.892 | 0.984 | 0.991 |
| NFI | 0.835 | 0.891 | 0.983 | 0.990 |
| NNFI | 0.770 | 0.840 | 0.971 | 0.983 |

# Model Comparison

| Models compared | df difference | Chi-square difference | Significance (p<.05) |
| --- | --- | --- | --- |
| One-factor vs. Correlated two-factor | 1 | 1104.92 | Significant |
| Correlated two-factor vs. Bi-factor | 3 | 1814.63 | Significant |
| Bi-factor vs. Correlated trait-uncorrelated method | 1 | 142.27 | Significant |

# Factor Loadings and Correlations

| Variables | Content | Manner | Long | Short |
|---|---|---|---|---|
| Grammar | 0.88* | | | |
| LongVoc | 0.73* | | 0.57* | |
| ShortVoc | 0.76* | | | |
| Pronunciation | | 0.83* | | 0.71* |
| ShortFluency | | 0.80* | | 0.18* |
| LongFluency | | 0.73* | 0.57* | |
| ToneReception | | 0.33* | | |
| ToneProduction | | 0.80* | | 0.43* |

| | Manner |
|---|---|
| Content | 0.915* |

*p<.05

ALWAYS LEARNING

PEARSON

# Research Questions

RQ 1: Does the SCT discriminate as expected according to known populations?

RQ 2: What is the factor structure of the SCT test scores?

RQ 3: Are SCT's subscores reliable enough to be taken as stable indicators of test-takers' strengths and weaknesses?

# *Test Reliability*

| Score | Split-half Method (n=166) | Test – Retest Method (n=158) | Human Scoring (n=166) |
|---|---|---|---|
| Overall | **0.97** | **0.95** | **0.98** |
| Grammar | 0.92 | 0.91 | 0.96 |
| Vocabulary | 0.94 | 0.93 | 0.96 |
| Fluency | 0.97 | 0.93 | 0.96 |
| Pronunciation | 0.96 | 0.91 | 0.95 |
| Tone | 0.93 | 0.87 | 0.96 |

# Summary

- SCT appears to measure facility in spoken Chinese, discriminating different test-taker groups (L1 Chinese, Heritage speakers, Non-heritage learners)

- SCT's MTMM is supported

- SCT's score reporting logic with content and manner aspects of language performance is supported

- SCT's reliability estimates are high and stable

# Implications

- Subscores do not provide any specific problem areas that could be useful for intervention

  e.g., no specific grammar points or phonetic points

- SCT subscores can still be good indicators of *relative* strengths and weaknesses

- Perhaps, conduct a survey with test score users to understand whether and how these subscores are useful for diagnostic purposes

# *Thank you*

**masanori.suzuki@pearson.com**
**yujie.jia@pearson.com**

PEARSON