# STANDARDIZING QUALITY ASSESSMENT FOR THE MULTILINGUAL WEB

Leonid Glazychev, Ph.D., CEO
Logrus International Corporation

# ASTM STANDARD PROPOSAL WK46397

- RATIONALE
  - Standards crucial for all stages of content production
    - Including quality assessment of multilingual materials
  - No methodology or metrics for public Language Quality Assurance (LQA)
  - Executive Order 13166: http://www.lep.gov/, http://www.justice.gov/crt/about/cor/13166.php
    "Improving Access to Services for Persons with Limited English Proficiency", reaffirmed in 2011
- WK46397
  - "Development of a complete methodology, including a simplified quality metric, for crowd-sourced expert language quality assessment targeted at nonprofit web sites and other documents of public interest."
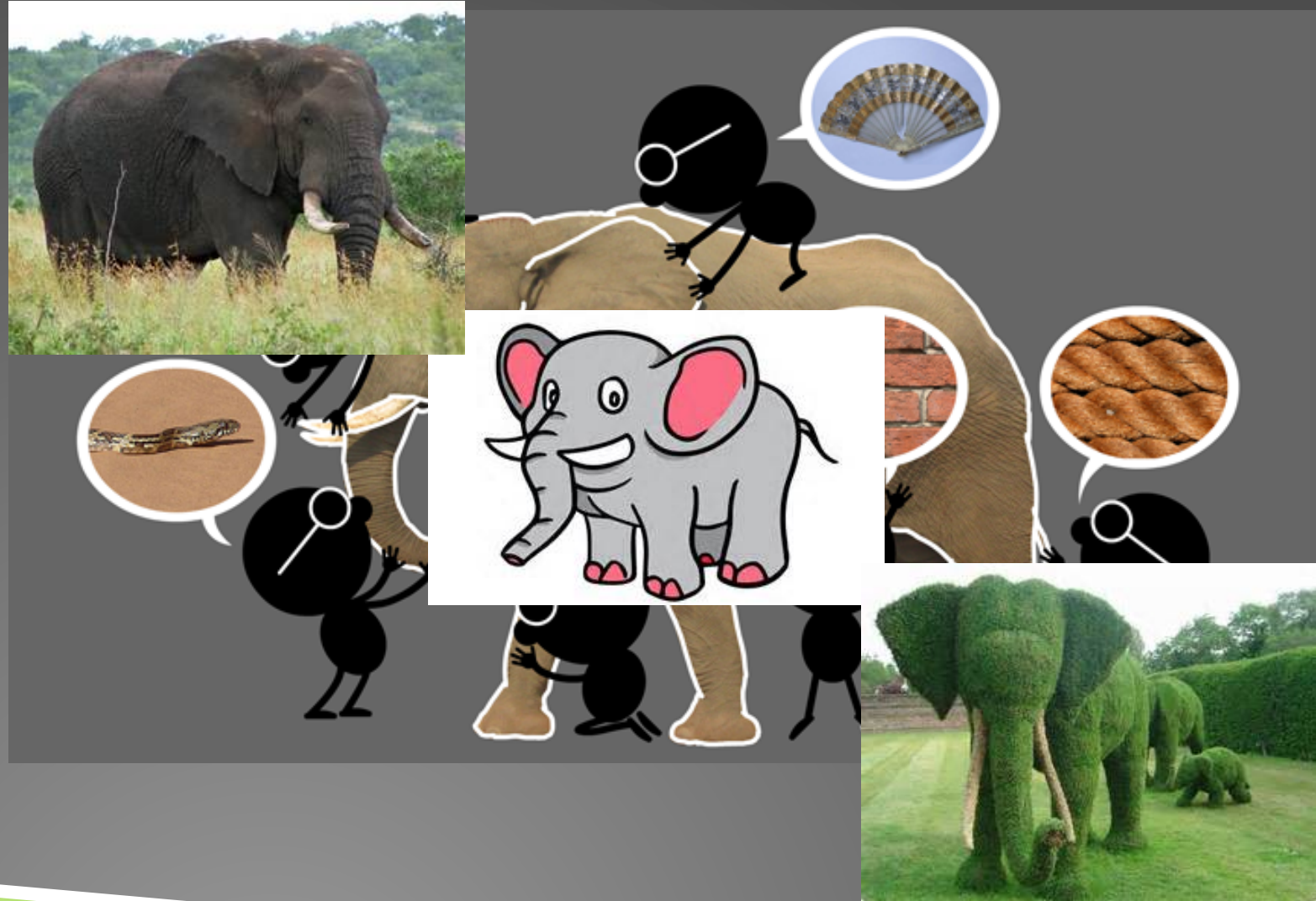- PRIMARY GOAL: A simplified quality assessment standard
  - Quick, inexpensive and reliable initial assessment
  - Reviewing big, highly visible resources designated for wide public use
    - Large and significantly diverse target audience
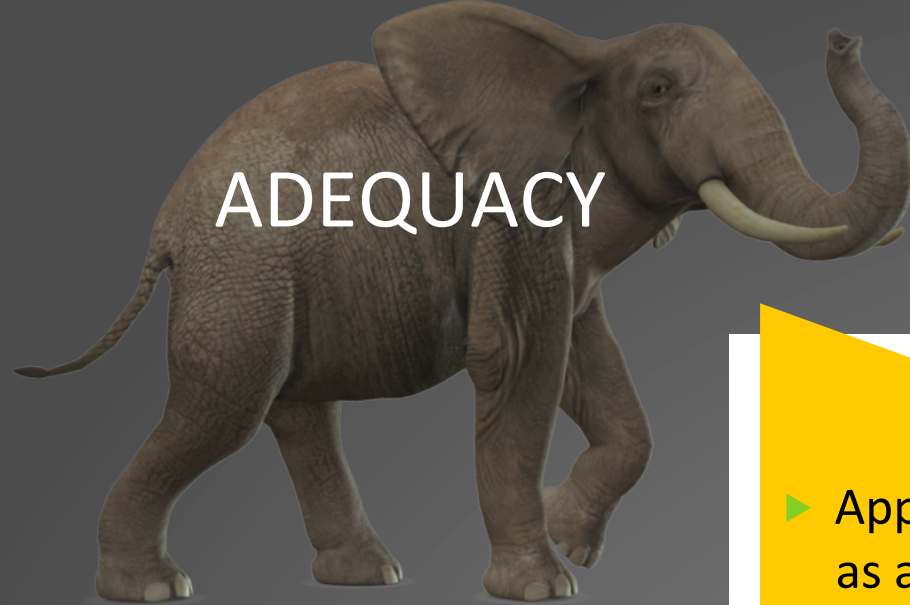    - Limited review capabilities and/or budget
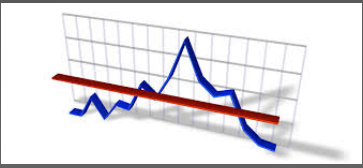
# EMPHASIS ON HOLISTIC ASSESSMENT

- The whole is always more important to us than its constituents

- Object properties can't be fully revealed or described based on its parts alone

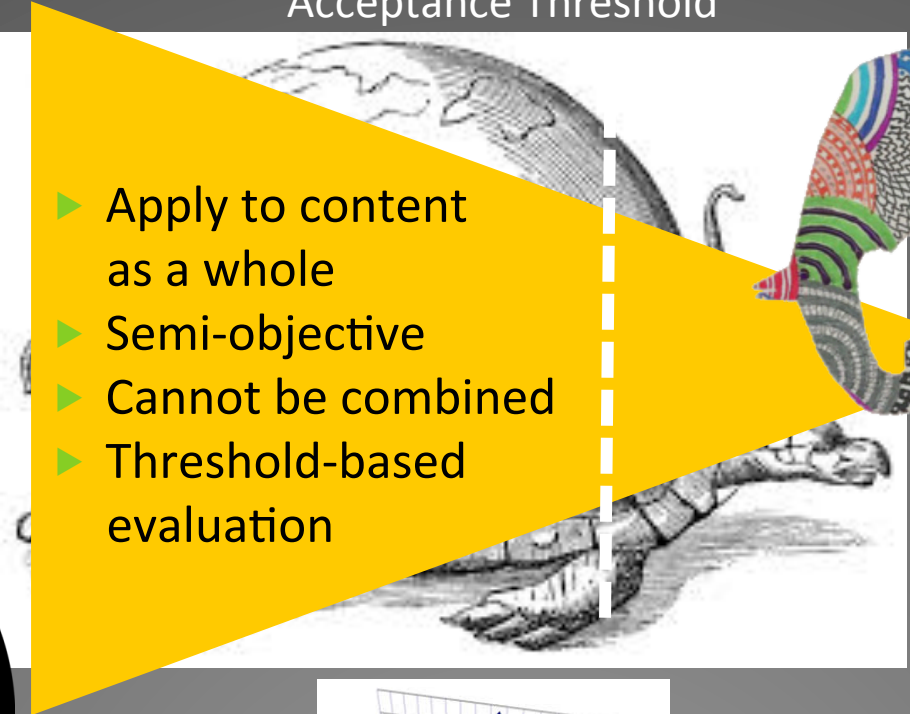- Quality assurance cannot be complete or accurate if there is no way of making holistic evaluations

# THE QUALITY TRIANGLE

Universal, 3D quality picture
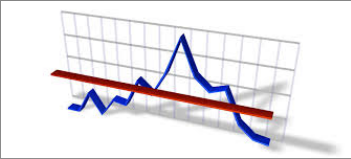- ▶ Same approach
- ▶ Any issue catalogue
- ▶ Only expectations vary

ADEQUACY

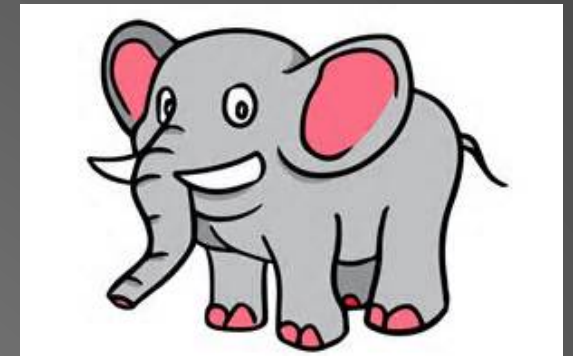Acceptance Threshold

HOLISTIC FACTORS

- ▶ Apply to content as a whole
- ▶ Semi-objective
- ▶ Cannot be combined
- ▶ Threshold-based evaluation

READABILITY

Acceptance Threshold

ATOMISTIC QUALITY

$$Q_A = \sum_{i=1}^{\infty} \frac{(W_i N_i)}{V}$$

# ATOMISTIC QUALITY

- Measured in ALL existing quality metrics
- Opposite to holistic
- Applies to
  - Quality issues at the "atomic" level of the content (vs. holistic)
  - Sentences, strings, translation units, …
- Includes issues like
  - Terminology inconsistency or deviations
  - Style guides, country standards
  - Tags, placeholders
  - Formatting
  - …
- Complements holistic usability/quality evaluation
- Example of a comprehensive issue framework
  - MQM: http://www.qt21.eu/launchpad/content/multidimensional-quality-metrics

# THE PRICE OF OBJECTIVITY

- Objective = Universal issue nature
  - Explanation doesn't require language knowledge
- No dependence on the reviewer's personality
  - A typo is still a typo
  - An error in country standards is still an error anyway
  - Everything depends on issue classification and the weighting system
- Price of objectivity
  - Comprehensive and clear issue classification
  - Availability of all ancillary materials
    - Glossaries, style guides, special requirements, etc.
  - Professional reviewers
  - Extensive language quality assurance (LQA) training
  - Detailed issue logging
  - Issue reconciliation with translators
  - Time and cost
  - **Otherwise the assessment is subjective and inaccurate!**
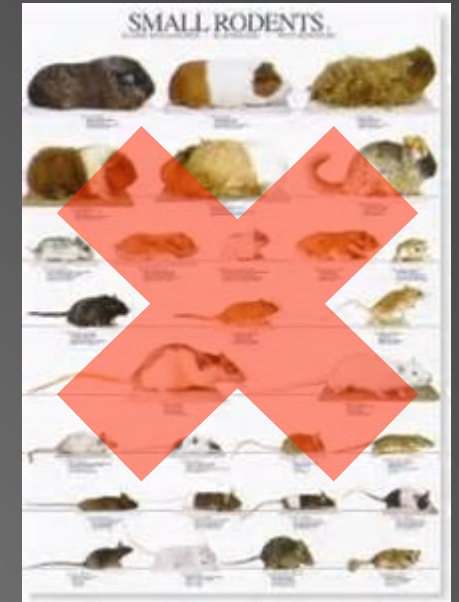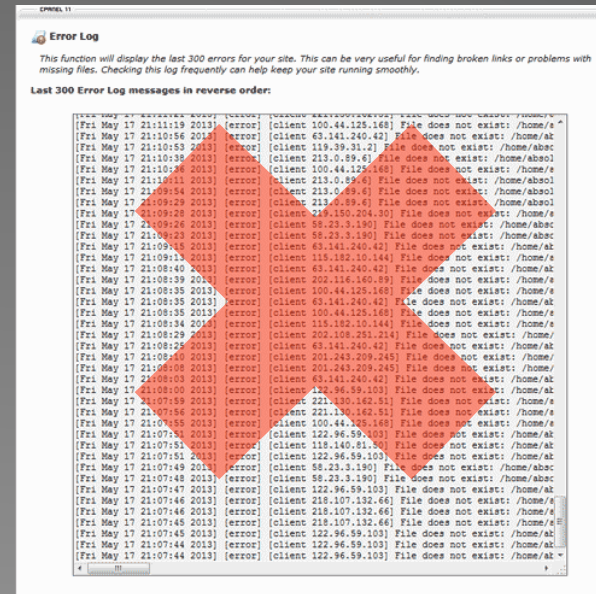
Wait, you want me to work 40 + hours a week UNPAID?

# THE LIMITATIONS

- Cannot expect serious preparation
  - Minimal/no reviewer training
  - Just explain the task in the simplest terms possible
- No thorough issue catalogues/quality frameworks
  - Unless completely trivial
- No serious quality issue logging
  - Just ask to provide typical examples
  - Make the feedback form simple and short

- **OUT OF THE QUESTION:**
  - **Complicated requirements**
  - **Strict definitions**
  - **Quality frameworks**
  - **Special rules, etc.**

▶ Defining all three cornerstones

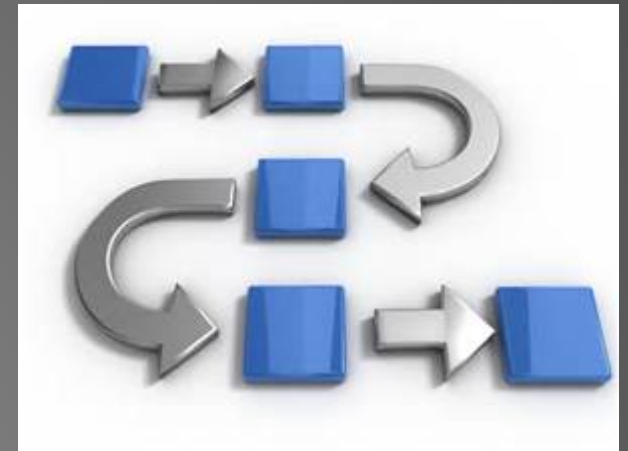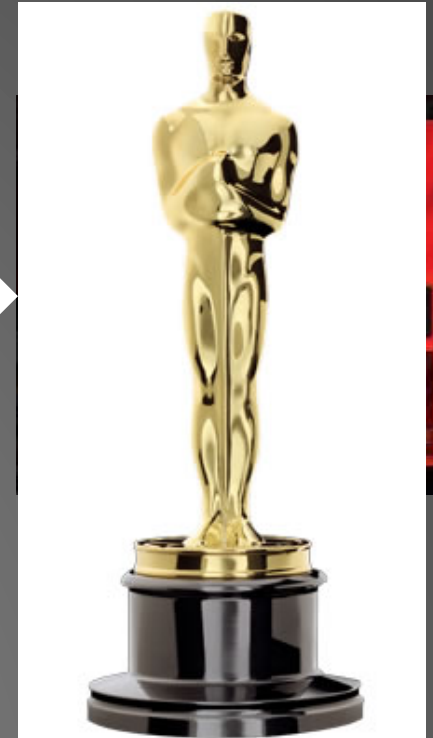General Approach/Methodology



Quality Metric



Process

# GENERAL APPROACH AND METHODOLOGY

- Simplified methodology
  - Focusing on holistic evaluations
- Objectivity and accuracy gained through statistics
  - Meaningful averages and standard deviations
  - Multiple people reviewing the same piece
  - Essential to collect sufficient statistics
- Limit contributors to language professionals only
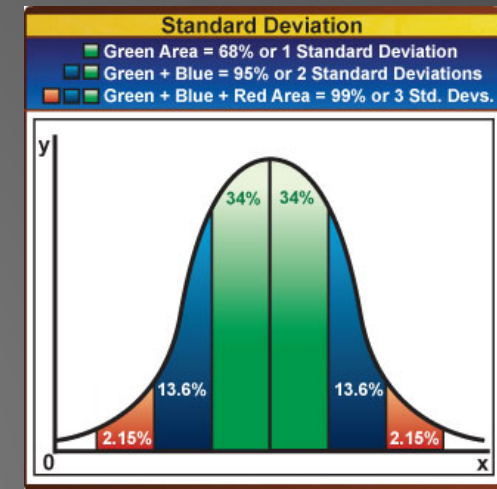
# SIMPLIFIED QUALITY SQUARE METRIC

- The Quality Square approach
- Simplified form, no detailed issue definitions or formal requirements
- Four ratings for each text on a 0-10 scale
  - The number of major (showstopper) errors
    - 0 => 10
    - 1 => 5
    - 2 or more => 0
  - Holistic translation readability
    - 0 = Completely unreadable/incomprehensible
    - 10 = Perfectly intelligible and readable text
  - Holistic translation adequacy
    - 0 = Completely inadequate
    - 10 = Perfectly conveyed meaning
  - Atomistic quality
    - 0 = Overabundance of atomistic-level errors
    - 10 = Completely error-free text
  - **A brief explanation required in each case**





PROHIBIDO EL PASO

404 error:
Building not found

# QUALITY ASSURANCE PROCESS

- Clear and brief LQA review scope
- Translated content frozen
- Online portal with project description/scope definition
- Pre-process results
- Individual list of pre-processing checks for each project
- Calculating median ratings and standard deviations
- Comparing all ratings against pre-defined thresholds

# EXPECTATIONS

▶ Reliable, statistically sound high-level LQA results in the crowdsourcing environment

▶ Cannot serve as a valid replacement for professional LQAs

| Untrained translators or linguists | Specially trained professionals |
|---|---|
| Almost no formal criteria | Extensive and well-defined formal criteria |
| General criteria | Criteria fine-tuned to the client's requirements |
| Minimal accuracy and consistency | High level of sophistication, accuracy, and consistency |

▶ Obtaining quick results at a minimal (or zero) cost

  ▶ Getting a rough evaluation of translation quality

  ▶ Reveal significant problems

  ▶ Assess the need for a professional LQA

▶ Acceptance thresholds replaced by "alarm-raising" ones

# CASE-STUDY: US ACA SPANISH WEBSITE REVIEW

- Originally requested directly by the US government
  - Affordable Care Act Spanish-Language Website: www.CuidadoDeSalud.gov
- Carried out free of charge by Logrus International for GALA
  - Globalization and Localization Association, www.gala-global.org
- Logrus developed and provided methodology
- Logrus organized the review and provided analytics
- Volunteer effort, crowdsourcing-based approach

# PROCESS ORGANIZATION

- Strictly following the process described earlier
- Simplified Quality Square methodology applied
  - Major errors (10 = None, 0 = More than 2)
  - Readability (0 - 10)
  - Adequacy (0 - 10)
  - Atomistic (0 – 10)
- 18 contributors chosen among language professionals only
- Mini-portal for participants
  - Self-registration
  - Brief error category definitions
  - Entering ratings and comments
- Comprehensive data pre-processing, discarding:
  - Standalone "perfect" (10 out of 10) evaluations
  - Marginally high or low ratings with no explanations
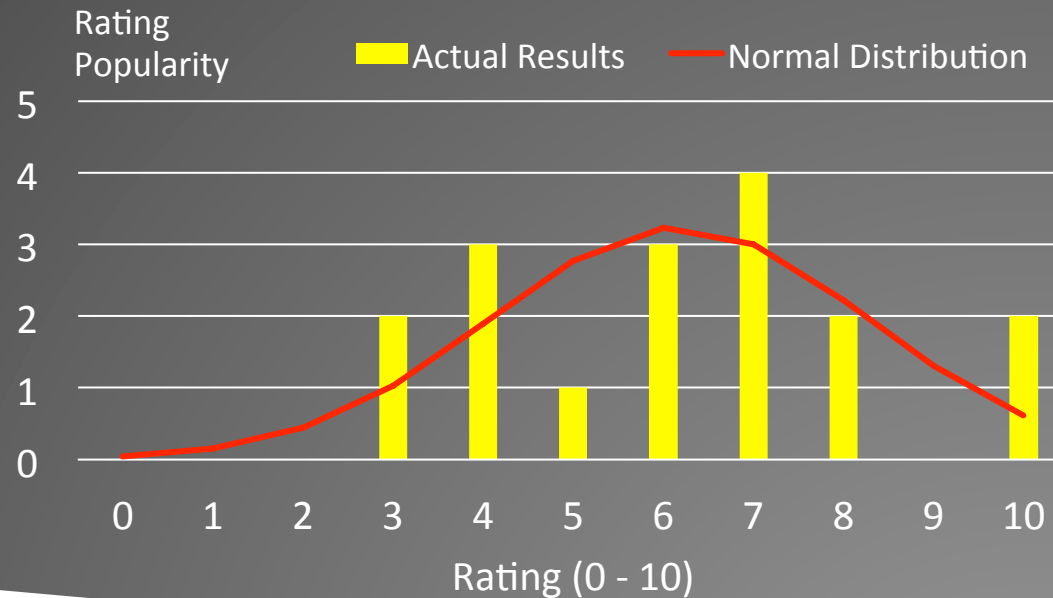  - Skewed ratings caused by reviewer errors

# PROJECT SPECIFICS

- Target language specifics
  - Most translation and LQA tasks target a specific region
    - Latin America (LatAm), Argentina, Mexico, Spain…
  - Each reviewer had a particular language "flavor" in mind
  - Target audience = Spanish-speaking population in the US
    - People with various backgrounds
    - Speaking a wide variety of Spanish, or even "Spanglish"
  - Most neutral and universal translation not sounding natural to some native speakers
- Understanding the review scope
  - Some "major errors" were functional issues beyond the LQA scope
    - Navigating health insurance plans and prices in English
    - Spelling errors in responses obtained through the chat feature
- Disregarded during pre-processing
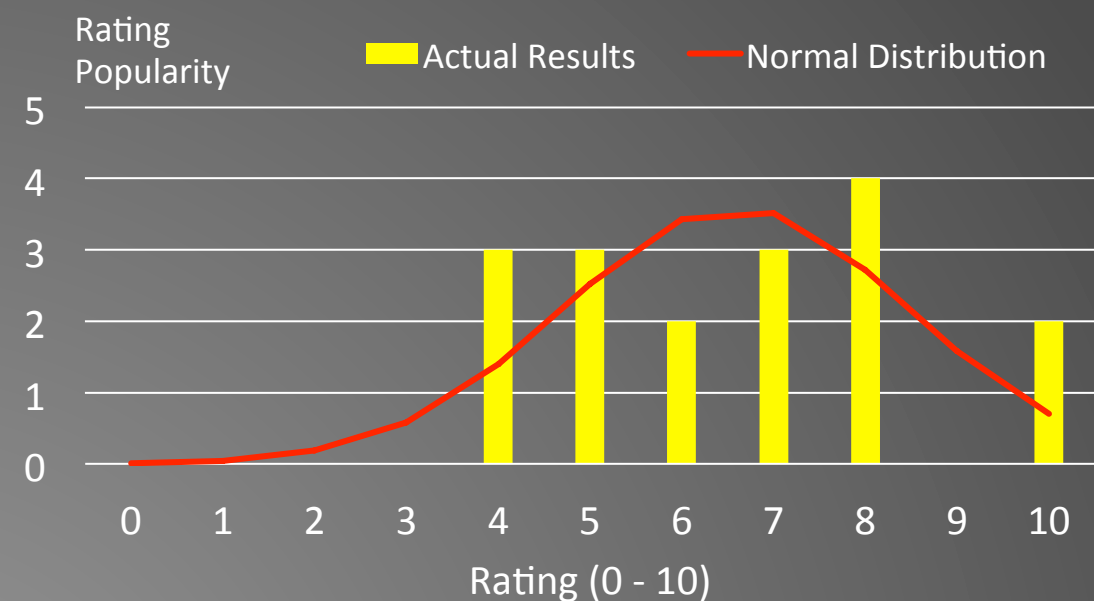  - Targeting translation quality alone, not portal usability or functionality

# MAJOR ERRORS – READABILITY – ADEQUACY

- Major errors: None (11), More than 2 (7), 1 grade ignored
- Readability and Adequacy
  - **YOUR** reviewer could contribute to ANY of the bars
  - Only threshold-based criteria really work

**Readability. Mean value: 6.2, Std. Deviation: 2.1**
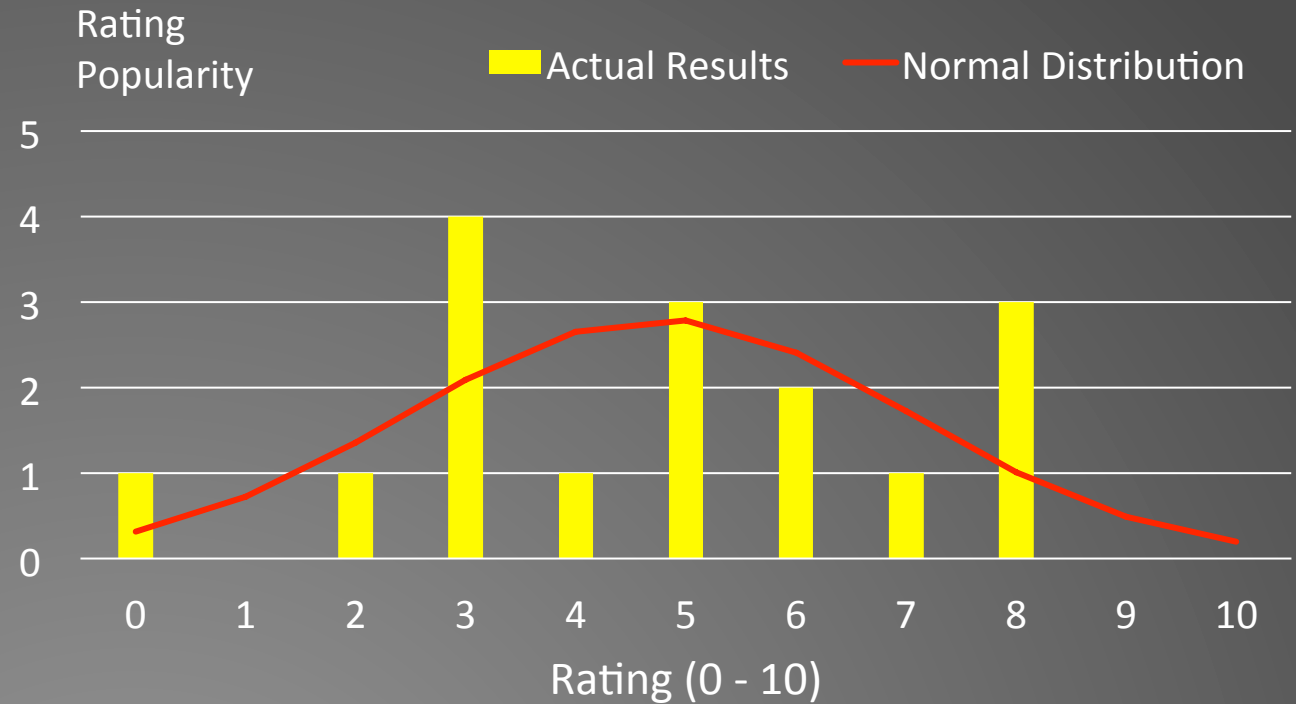
**Adequacy. Mean value: 6.6, Std. Deviation: 1.9**

# ATOMISTIC QUALITY

- Biggest opinion spread
  - Illustrates the gap between professional and crowd-sourced work
  - No detailed criteria or training
  - Should be the most objective factor ☺
- "Mechanical" stats
  - Mean value: 5.4
  - Standard deviation: 2.8
- Adjusted stats
  - Mean value: 4.7
  - Standard deviation: 2.4

Atomistic Quality.  Mean value: 4.7, Std. Deviation: 2.4

# METRIC AND PROCESS SUMMARY

- Both holistic Readability and Adequacy results can be relied upon
- Good basis for assessing overall translation quality
- Judgment about the presence of Showstopper errors is convincing
- Atomistic quality assessment is not accurate enough
  - Gives a good general idea of the pervasiveness of non-critical, atomistic-level errors

- Major crowdsourcing LQA results look trustworthy and consistent
- A reliable high-level picture of translation quality
- Experimental proof that the whole model works
  - Even in the relatively extreme crowdsourcing environment