



UNIVERSIDAD DE  
COSTA RICA

ELM

Escuela de  
Lenguas Modernas

PELEx

Programa de  
Evaluación en  
Lengua Extranjera

ALTE

Associate Member

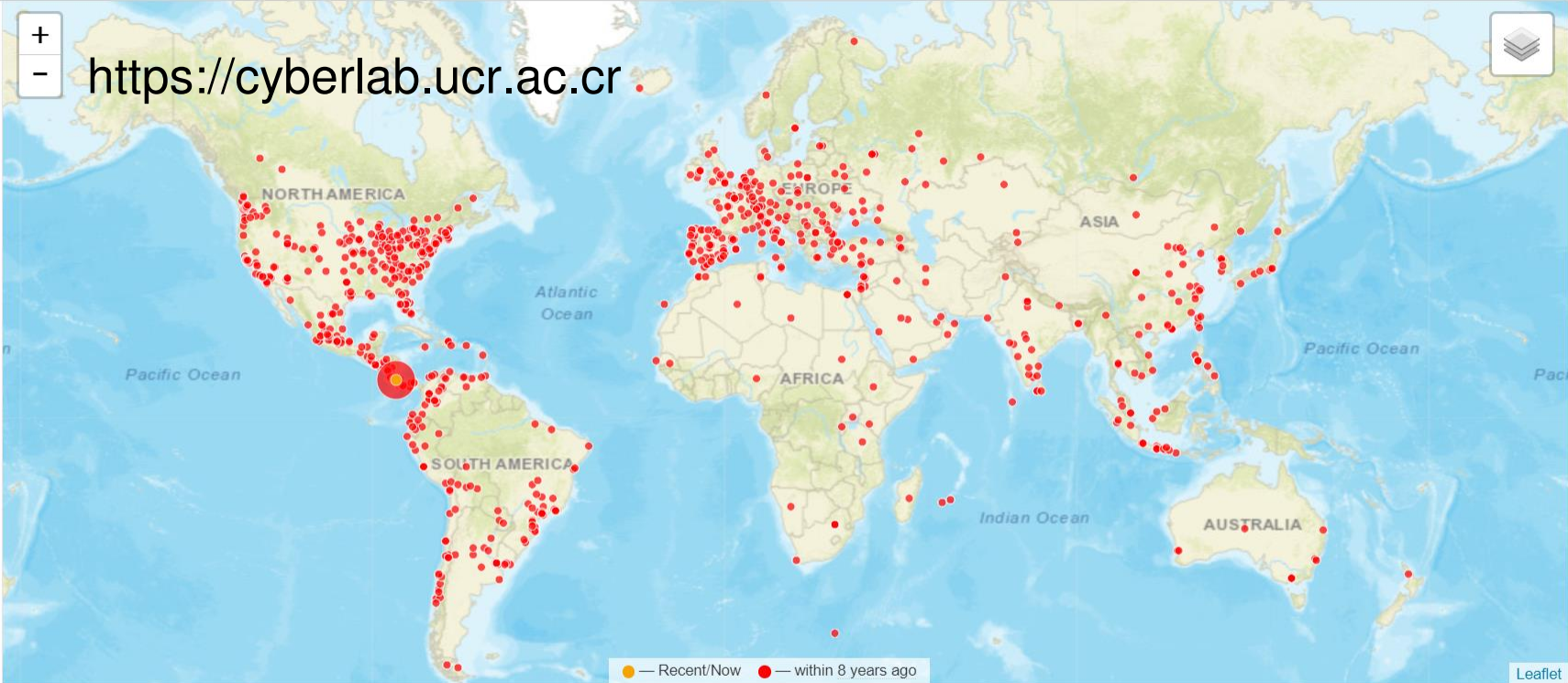
# Implicaciones y desafíos del uso de Inteligencia Artificial en la evaluación por competencias en lenguas extranjeras

EL CASO DE COSTA RICA



Allen Quesada Pacheco, Ph.D

2023



**Recent visitors** **Browser** **OS** **Date**

 <b>New</b> Visitor from San José, Costa Rica with 1 pageview	 Chrome	 Win10	December 5, 2023, 12:47 pm
--	--	---	----------------------------

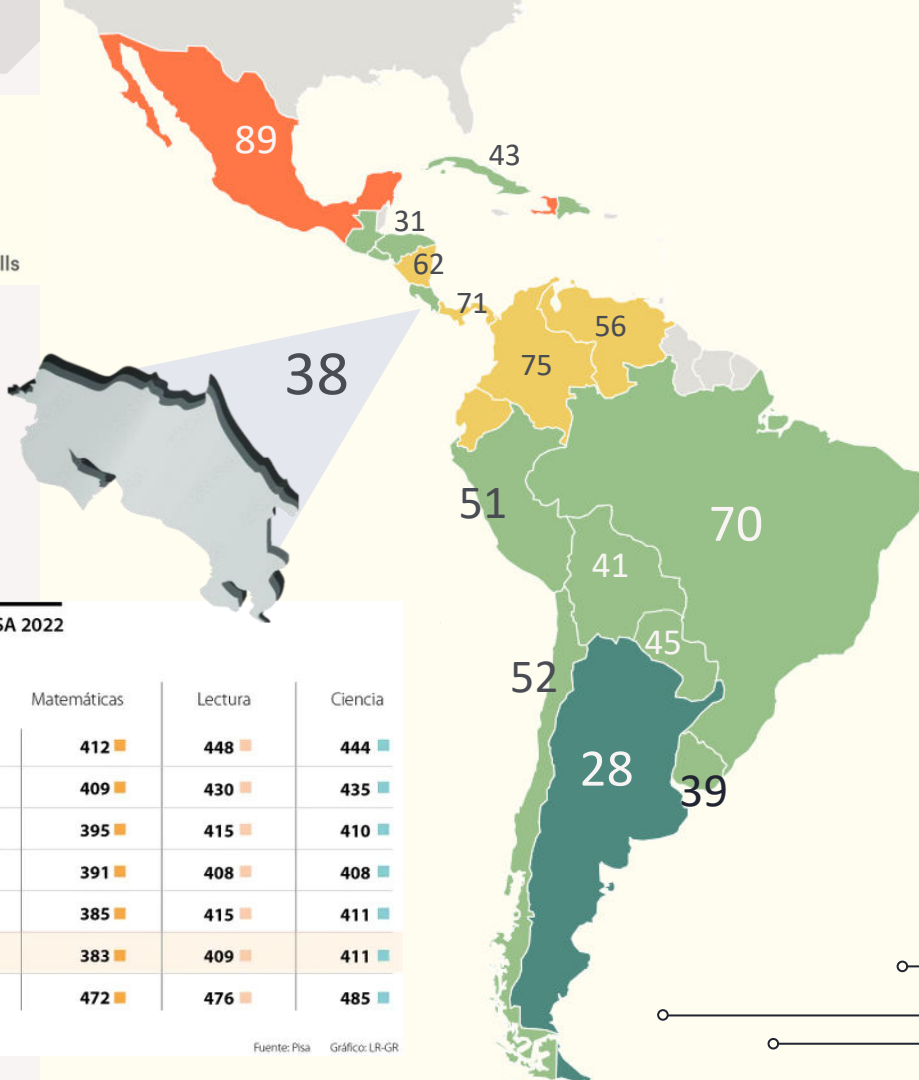


# EF EPI

## EF English Proficiency Index

A Ranking of 111 Countries and Regions by English Skills

- Very high
- High
- Moderate
- Low
- Very low

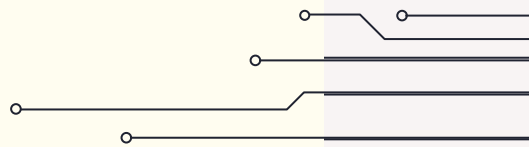


### RESULTADOS PISA 2022

#### AMÉRICA LATINA

	Matemáticas	Lectura	Ciencia
1 Chile	412	448	444
2 Uruguay	409	430	435
3 México	395	415	410
4 Perú	391	408	408
5 Costa Rica	385	415	411
6 Colombia	383	409	411
Promedio Ocde	472	476	485

Fuente: Pisa Gráfico: LR-GR





# EF EPI

## EF English Proficiency Index

A Ranking of 111 Countries and Regions by English Skills



525 543

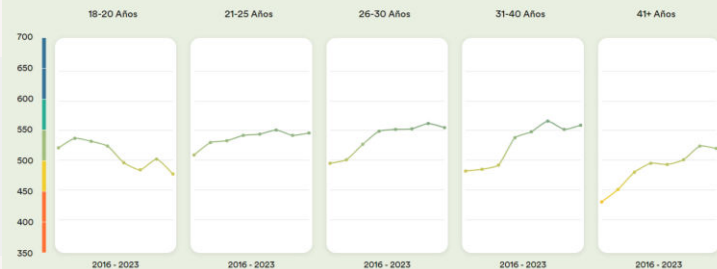
Región principal: Heredia\* (566)

Ciudad principal: Heredia (565)

Grupo de edad principal: 31-40 (nota media de 559)



### Tendencias según edad



- Very high
- High
- Moderate
- Low
- Very low

### Costa Rica

# #38

 de 113

Clasificación EF EPI: 534

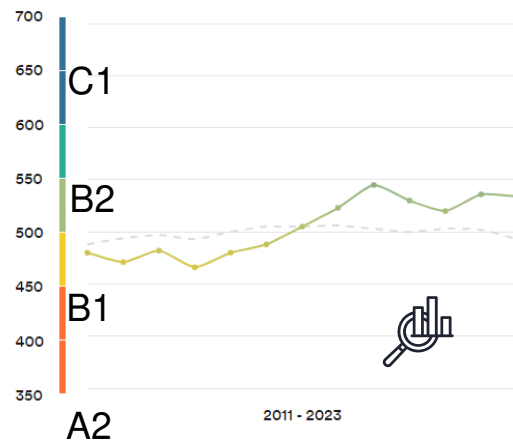
Puntuación media global: 493

Posición en Latinoamérica: 3 de 20

[Descargar hoja de datos](#)



### Tendencias del EF EPI

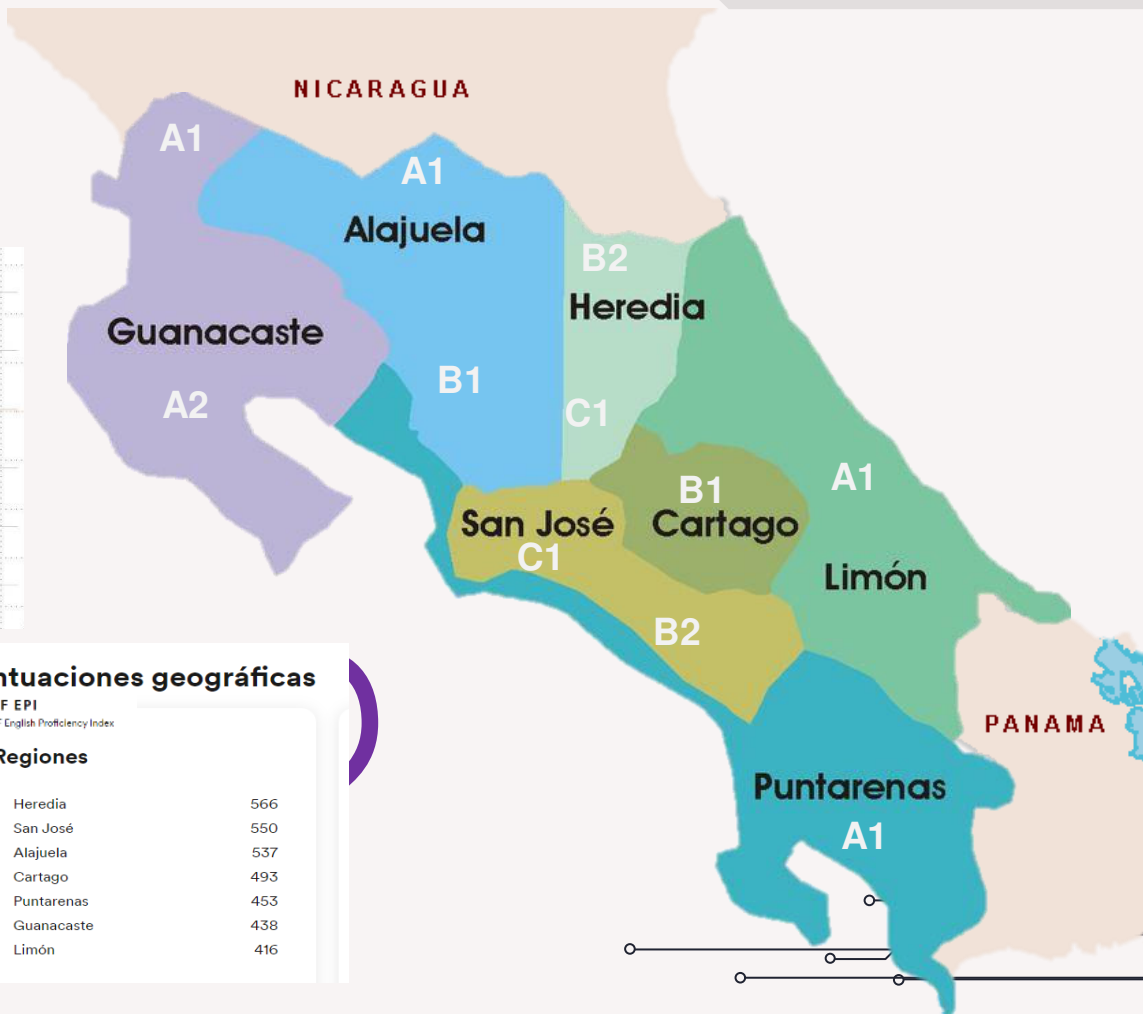
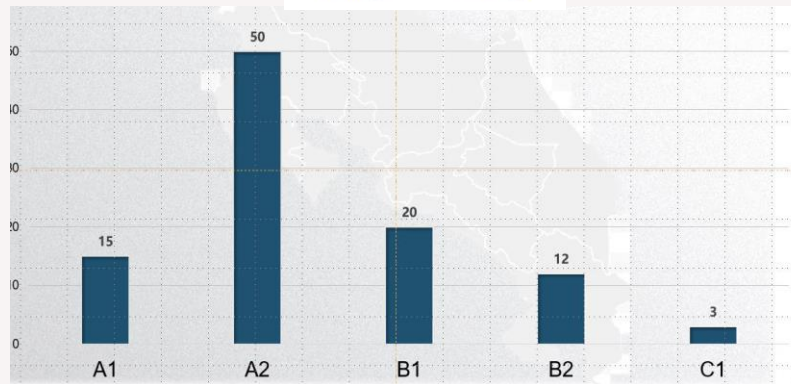




Reading



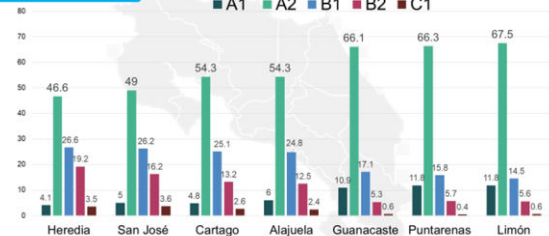
Listening



**PDL 2021**

Resultados porcentuales por provincia

■ A1 ■ A2 ■ B1 ■ B2 ■ C1



Fuente: Prueba de dominio lingüístico PDL 2021, Escuela de Lenguas Modernas UCR

**Puntuaciones geográficas**

EF EPI  
EF English Proficiency Index

Regiones

Heredia	566
San José	550
Alajuela	537
Cartago	493
Puntarenas	453
Guanacaste	438
Limón	416

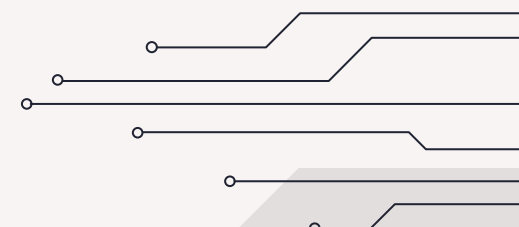
# Marcos de Referencias dominio de idiomas

El Marco Común Europeo de Referencia (MCER) es un marco de referencia de reconocimiento internacional para describir el dominio de un idioma (A1-C1).



## Hay varios marcos de referencia con objetivos similares:

- Consejo Americano para la Enseñanza de Lenguas Extranjeras (ACTFL, American Council on the Teaching of Foreign Languages Proficiency Guidelines),
- Los exámenes de competencias lingüísticas canadienses (CLB, Canadian Language Benchmarks),
- La escala de competencia de la ILR (Interagency Language Roundtable) y
- WIDA (Wisconsin, Delaware, Arkansas Assessment language framework)



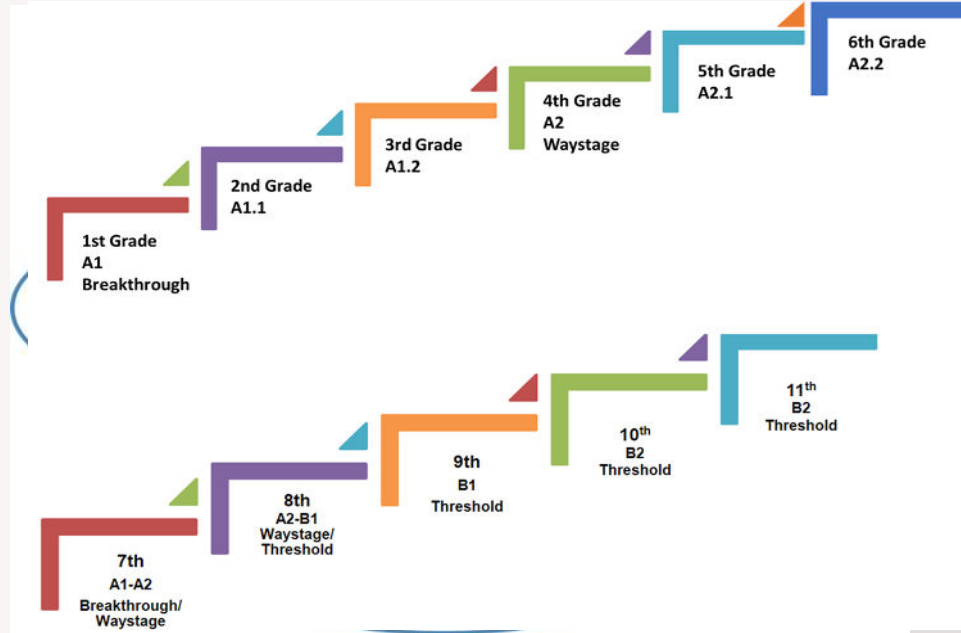


UNIVERSIDAD DE COSTA RICA

mep  
Ministerio de Educación Pública



speechace







**A L T E**

Associate Member  
Adaptativa y Automatizada con IA

Digital Evolution Education Program  
Allen Quesada Pacheco, PhD  
Universidad de Costa Rica - 2021

**PURA VIDA**

Laspau




**ECCI**

---

Escuela de  
**Ciencias de la  
Computación e  
Informática**



**Laspau** Affiliated with  
Harvard University

**LAALTA**  
Latin American Association  
for Language Testing and Assessment



# Rubrics

## Criteria Oral Evaluation Sheet LM-1002 / Quiz & Exam / B1 Intermediate-mid level in academic contexts

Content 20% Grammar 20% Vocabulary 20%  
Pronunciation 20% Fluency 10% Communication skills 10%

Each aspect in the evaluation has its own percentage; however, poor performance in any of the aspects will affect the others (e.g. poor pronunciation affects comprehensibility of content, poor

Score	Content
Excellent (100-90) 20-16pts	Can produce with some confidence <b>connected ideas</b> to give <b>detailed accounts</b> of experiences, events and surrounding circumstances about <b>all</b> the topics studied in the course. Can briefly give <b>reasons</b> and explanations to <b>justify</b> opinions and plans on the unit topics. Can report straightforward information to indicate the nature of a problem, and explain the main points with <b>reasonable precision</b> . Can <b>ask for further details</b> on the topic. Can take part in formal <b>discussion</b> of familiar subjects exchanging <b>plenty</b> information and discussing solutions to problems.
Good (89-80) 17,8-16pts	Can produce with some confidence <b>connected ideas</b> to give <b>detailed accounts</b> of experiences, events and surrounding circumstances about <b>most</b> of the topics studied in the course. Can briefly give <b>reasons</b> and explanations to <b>justify most</b> opinions and plans on the unit topics. Can report <b>straightforward information</b> to indicate the nature of a problem <b>most of the times</b> , and explain the main points with <b>reasonable precision</b> . Can <b>ask for further details in most instances</b> . Can take part in formal <b>discussion</b> of most familiar subjects exchanging <b>enough</b> information and <b>sometimes</b> discussing <b>solutions</b> to problems.
Acceptable (79-70) 15,8-14pts	Can produce with some confidence <b>slightly connected ideas</b> to give <b>some detailed accounts</b> of experiences, events and surrounding circumstances and explanations to indicate the topics studied in the course. Can briefly give <b>reasons</b> and explanations to <b>indicate the nature of a problem</b> <b>some of the times</b> , and explain the main points with <b>some</b> information. Can <b>ask for further details in some instances</b> . Can take part in formal <b>discussion</b> of familiar subjects exchanging <b>some</b> information and <b>rarely</b> discussing solutions to problems.
Weak (69-55)	Can produce with <b>little confidence</b> <b>somehat</b> accounts of experiences, events and surrounding circumstances.

Name: \_\_\_\_\_  
 Movie: \_\_\_\_\_  
 Topic: \_\_\_\_\_  
 Index card quiz ( ) Impromptu speech ( ) Prepared Speech ( ) Exam I ( ) Exam II ( )  
 Grade: \_\_\_\_\_

Content 50%	Language 50%		
	Pronunciation 20%	Grammar 20%	Vocabulary 10%
5 (25%)	5 (10%)	5 (10%)	5 (5%)
6 (30%)	6 (12%)	6 (12%)	6 (6%)
7 (35%)	7 (14%)	7 (14%)	7 (7%)
8 (40%)	8 (16%)	8 (16%)	8 (8%)
9 (45%)	9 (18%)	9 (18%)	9 (9%)
10 (50%)	10 (20%)	10 (20%)	10 (10%)

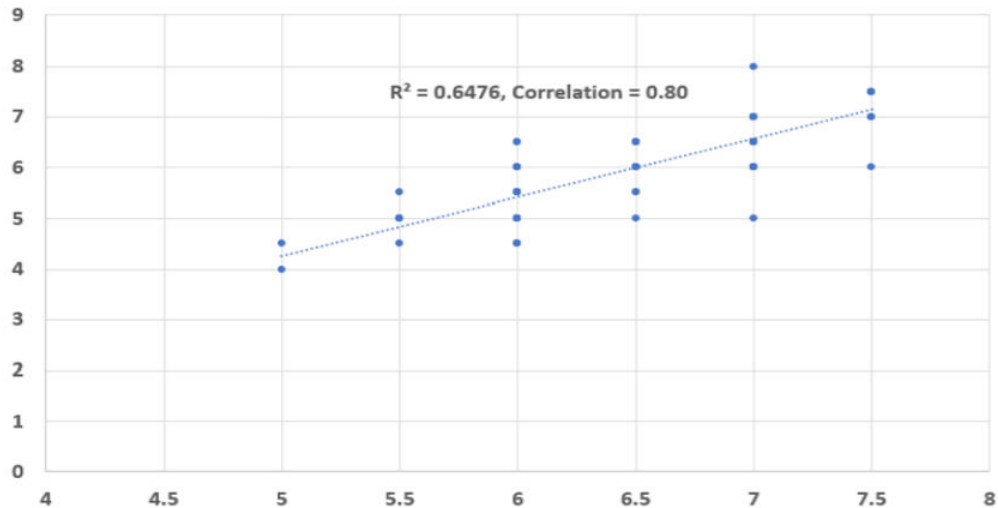
Errors in vocabulary or lack of it will affect content. / Errors in grammar and pronunciation will affect content. / Poor or inappropriate content will affect all the other areas.

Grammar & Vocabulary

Pronunciation

Skill components	Descriptors	Comments
<b>Pronunciation / Intonation</b>		
0.5 1 1.5 2 2.5 low High	<ul style="list-style-type: none"> <li>— Always</li> <li>— Most of the time</li> <li>— Much of the time</li> <li>— Sometimes</li> <li>— Rarely</li> </ul>	
<b>Fluency</b>		
0.5 1 1.5 2 2.5 low High	<ul style="list-style-type: none"> <li>— Uses native-like flow of speech</li> <li>— Uses fluent connected speech</li> <li>— Uses fluent connected speech, occasionally disrupted by search for correct form of expression</li> <li>— Speech is connected but frequently disrupted by search for correct form of expression</li> <li>— Uses simple sentences</li> </ul>	
<b>Vocabulary</b>		
0.5 1 1.5 2 2.5 low High	<ul style="list-style-type: none"> <li>— Uses sophisticated vocabulary in a variety of contexts</li> <li>— Uses varied and descriptive language, possibly including native-like phrasing and / or idiomatic expressions</li> <li>— Uses vocabulary sufficient to communicate in most social and academic contexts</li> <li>— Uses vocabulary sufficient to express needs and feelings and responds in familiar contexts</li> <li>— Uses only basic vocabulary with possible use of first language</li> </ul>	
<b>Grammar</b>		
0.5 1 1.5 2 2.5 low High	<ul style="list-style-type: none"> <li>— Uses appropriate tenses, pronouns, gender and number agreement, negation, articles, prepositions, and adjective placement</li> <li>— Always</li> <li>— Most of the time</li> <li>— Much of the time</li> <li>— Sometimes</li> <li>— Rarely</li> </ul>	

Teacher score vs Speechace score



speechace

UCR



Scores within  
 **$\pm 0.5$  IELTS**  
points of  
human grader



**0.8**  
Correlation with  
human grader

For Windows  
Ve a Configuración para activar Window

# TEYL/CAT Overview

1409



51%



8050



49%



### Table

Element	Field	Description
syllable_score_list[]	stress_level	Expected (Correct) lexical stress level based on the Lexicon. <b>Values:</b> <ul style="list-style-type: none"> <li>0: unstressed</li> <li>1: primary stress</li> <li>2: secondary stress</li> </ul>
syllable_score_list[]	predicted_stress_level	Detected lexical stress level based on the user audio. <b>Values:</b> <ul style="list-style-type: none"> <li>0: unstressed</li> <li>1: primary stress</li> <li>2: secondary stress</li> </ul>
syllable_score_list[]	stress_score	Floating point number between 0 and 100 indicating the correctness of the user's stress level.
syllable_score_list[]	pitch_range[]	[begin_pitch,end_pitch] recorded for this syllable in Hertz.
word_intonation_list[]	syllable_intonation_list[]	[pitch_change_from_previous, pitch_change_in_current] <b>pitch_change_from_previous:</b> indicates whether the pitch of a word rises, falls, or more or less unchanged. If we can't recognize the syllable (due to error from our side or due to user not saying the syllable), the value is null. If we can recognize the syllable but couldn't determine the pitch (due to error from our side or due to user reducing the sound to unvoiced), the value is REDUCED. <b>pitch_change_in_current:</b> is null unless one of the following cases occurs: If the pitch of the syllable falls, but the starting pitch of the syllable is higher than the ending pitch of previous syllable, secondary intonation is RISE while primary intonation is FALL. If the pitch of the syllable rises, but the starting pitch of the syllable is lower than the ending pitch of the previous syllable, secondary intonation is FALL while primary intonation is RISE.

### Table

Field	Description
duration	total length of speech
articulation	total length of articulation and non-speech events silence on very first segment.
speech_rate	speaking rate in syllables
articulation_rate	articulation rate in syllables
syllable_count	Count of syllables in this segment.
word_count	Count of words in this segment.
correct_syllable_count	Count of correctly spoken syllables
correct_word_count	Count of correctly spoken words
syllable_correct_per_minute	correct_syllable_count / duration
word_correct_per_minute	correct_word_count / duration
all_pause_count	count of all pauses (filler words) exceeding minimum pause threshold
all_pause_duration	total duration of all pauses
all_pause_list[]	a list of all the pauses exceeding 10 msecs
mean_length_run	mean length of run in syllables
max_length_run	max length of run in syllables
segment_metrics_list[]	A list of segments with metrics for each segment



## Test design

No.	Question Type	Weight	Max score	Rubric
1	Read aloud	10%	B2	50% Pronunciation 50% Fluency
2	Task achievement Ask a question	15%	B2	50% Pronunciation & Fluency 50% Task achievement
3	Task achievement Describe 4 picture story	20%	B2	50% Pronunciation & Fluency 50% Task achievement
4	Task achievement Share an opinion and rationale	25%	B2	50% Pronunciation & Fluency 50% Task achievement
5	Open ended	30%	B2	25% Pronunciation 25% Fluency 25% Vocabulary 25% Grammar



There are only two seasons in Costa Rica, the rainy season and the dry season. People usually enjoy going to see a volcano or they can go to the beach on sunny days. I like both seasons. We can have a lot of fun!

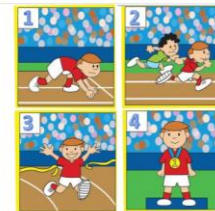
Read aloud the text on the right (Lea en voz alta)

Preparation timer 30



Ask one question about the picture. (Haga una pregunta de la imagen)

Preparation timer 30



Describe the story in the pictures. (Cuente una historia de las imagines)

Preparation timer 30



Give an opinion on the following topic: Do you believe recycling is important? Why or why not? (Que opinas de reciclar, y por que?)

Preparation timer 30



Talk about your future vacation. Where to go? Who with? How will you get there? (Hable sobre sus vacaciones en el futuro)

Preparation timer 30

### Fuentes Gonzalez Mariangel

Aug 23, 2023

[Link to secure online certificate](#)

[Speaking test report for Fuentes Gonzalez Mariangel \(Aug/23/2023\)](#)

Note: This link can be e-mailed and shared with others.



#### Summary of scores

Speechact **CEFR** IELTS PTE TOEFL TOEIC

CEFR score  
**A2**



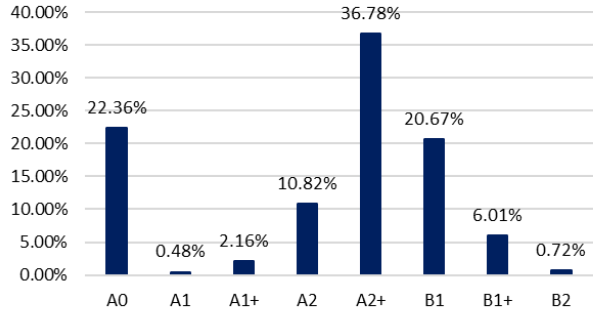
#### Descriptive feedback

Overall (A2)

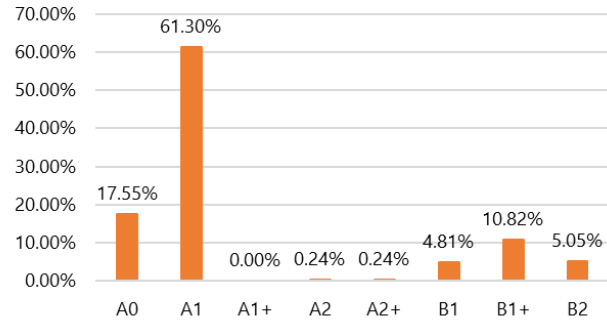
Demonstrates below average fluency and coherence. May need remedial training in speaking fluently and language construction. Has limited vocabulary and and may have difficulty expressing complex thoughts. Moderate task achievement but with significant errors. Has below average pronunciation accuracy and may not be easy to understand. Good at speaking simple sentences but regularly makes grammatical mistakes.

Pronunciation (C1)

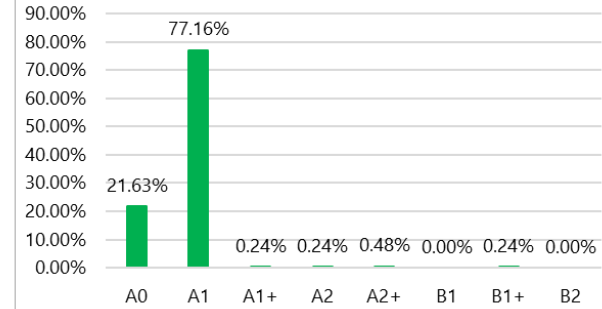
### Q1 – Read Aloud



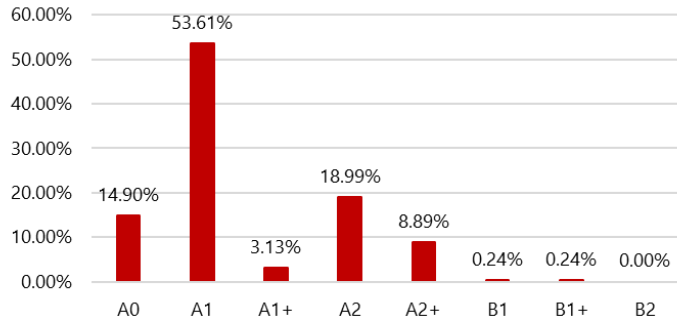
### Q2 – Ask a question



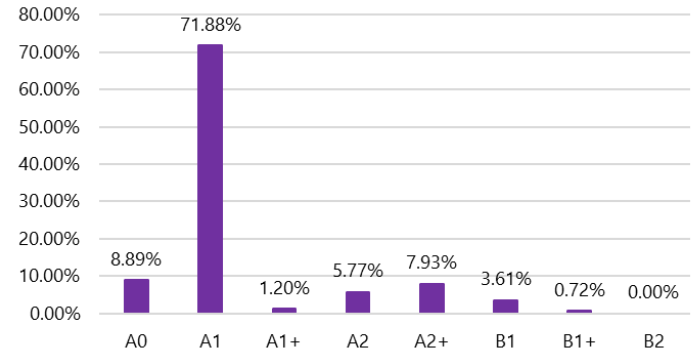
### Q3 – Describe 4 pic story



### Q4 – Share an opinion



### Q5 – Open ended

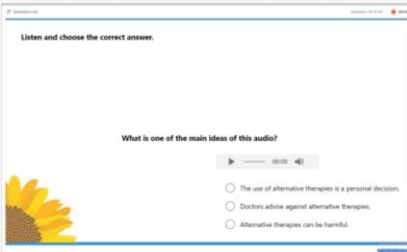
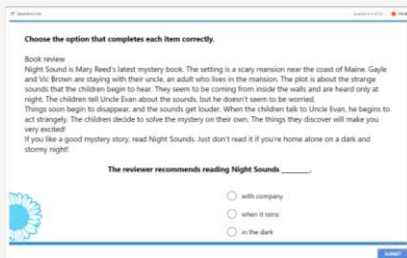




# Prueba Inglés Adaptativa

● Escucha

Lectura



# PeLexCAT: Adaptive Language Testing

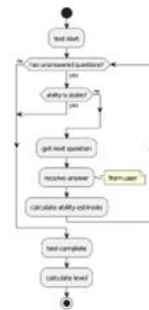
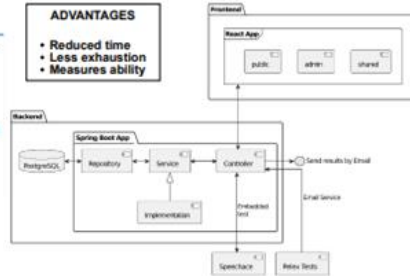
Development of a Localized Adaptive Language Test of Receptive Skills

Dr. Edgar Casasola Murillo and Dr. Allen Quesada Pacheco  
UNIVERSIDAD DE COSTA RICA



- SOFTWARE DESIGN**
- Scalable
  - Multi Layered
  - Email notifications
  - Accessibility
  - Integration of External AI Speaking Testing

- ADVANTAGES**
- Reduced time
  - Less exhaustion
  - Measures ability



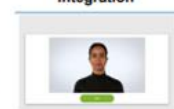
**ADAPTIVE TEST  
MAIN ALGORITHM**

- STOP CONDITION**
1. Fixed threshold
  2. Standard Error of Estimation <= 0.35
  3. Time limit

READING

MULTI SKILL TESTING  
INTEGRATION

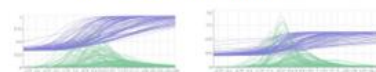
External System  
Integration



ORAL SPEAKING  
IA TESTING

LISTENING

3 Parameter IRT Model  
Discrimination vrs Difficulty



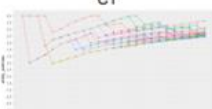
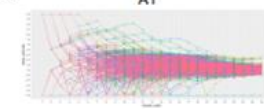
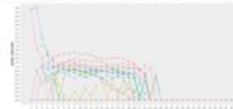
READING				
A1	A2	B1	B2	C1
30	27	28	20	21

LISTENING				
A1	A2	B1	B2	C1
27	21	26	17	15

**SIMULATION BASED EXPERIMENTATION AND VALIDATION**

INFORMATION BASED ITEM DISPATCHER

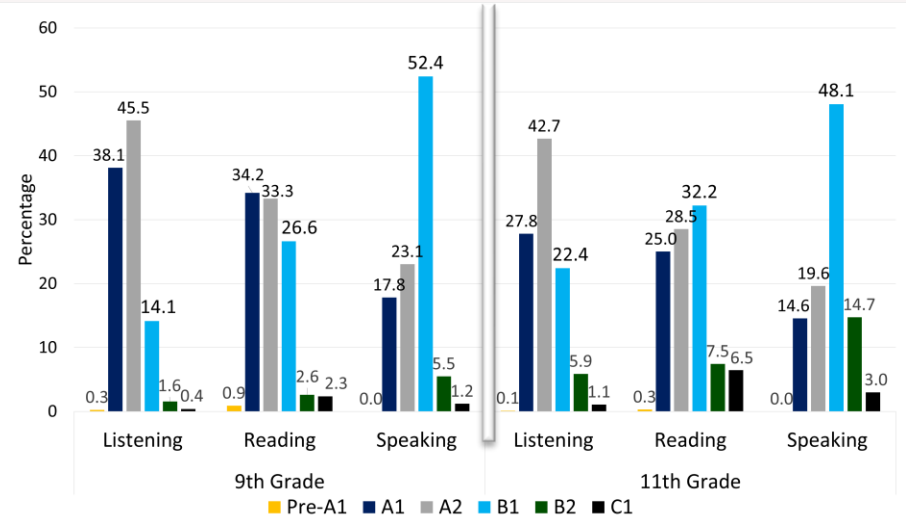
ABILITY CONVERGENCE BY LEVEL  
A1 C1



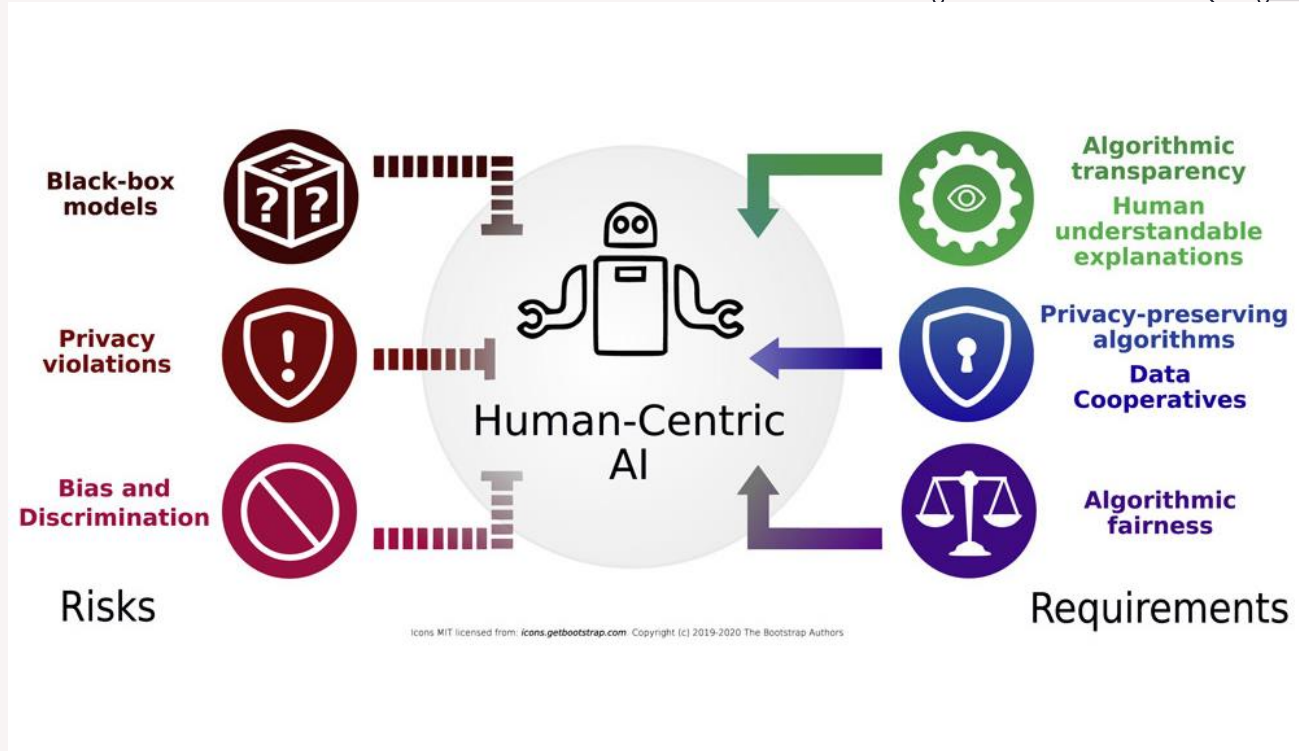
# TEYL

Banda	Reading	Listening	Speaking
Pre-A1	9.27	12.98	10.90
A1	56.94	58.31	60.19
A1+	16.58	18.57	13.80
A2	17.20	10.14	15.20
<b>Total</b>	100	100	100

# Secondary



# Ética en la IA



Lepri, B., Oliver, N., & Pentland, A. (2021). Ethical machines: The human-centric use of artificial intelligence. *iScience*, 24(3), 102249. doi:10.1016/j.isci.2021.102249

# LA INTELIGENCIA GENERATIVA (IG) EN LA EVALUACIÓN DE IDIOMAS

La IG puede aportar los siguientes beneficios a la evaluación de idiomas:

- Personalización:** La IG puede utilizarse para crear evaluaciones personalizadas para cada estudiante.

- Eficiencia:** La IG puede automatizar muchas tareas relacionadas con la evaluación, lo que puede liberar tiempo a los profesores para centrarse en otras actividades.

La IG se puede utilizar para aplicaciones específicas en la evaluación de idiomas, como:

- Generación de preguntas:** La IG se puede utilizar para generar preguntas de evaluación, lo que puede ayudar a crear evaluaciones más desafiantes y relevantes.

- Corrección de exámenes:** La IG se puede utilizar para corregir exámenes, lo que puede liberar tiempo a los profesores para centrarse en otras tareas.

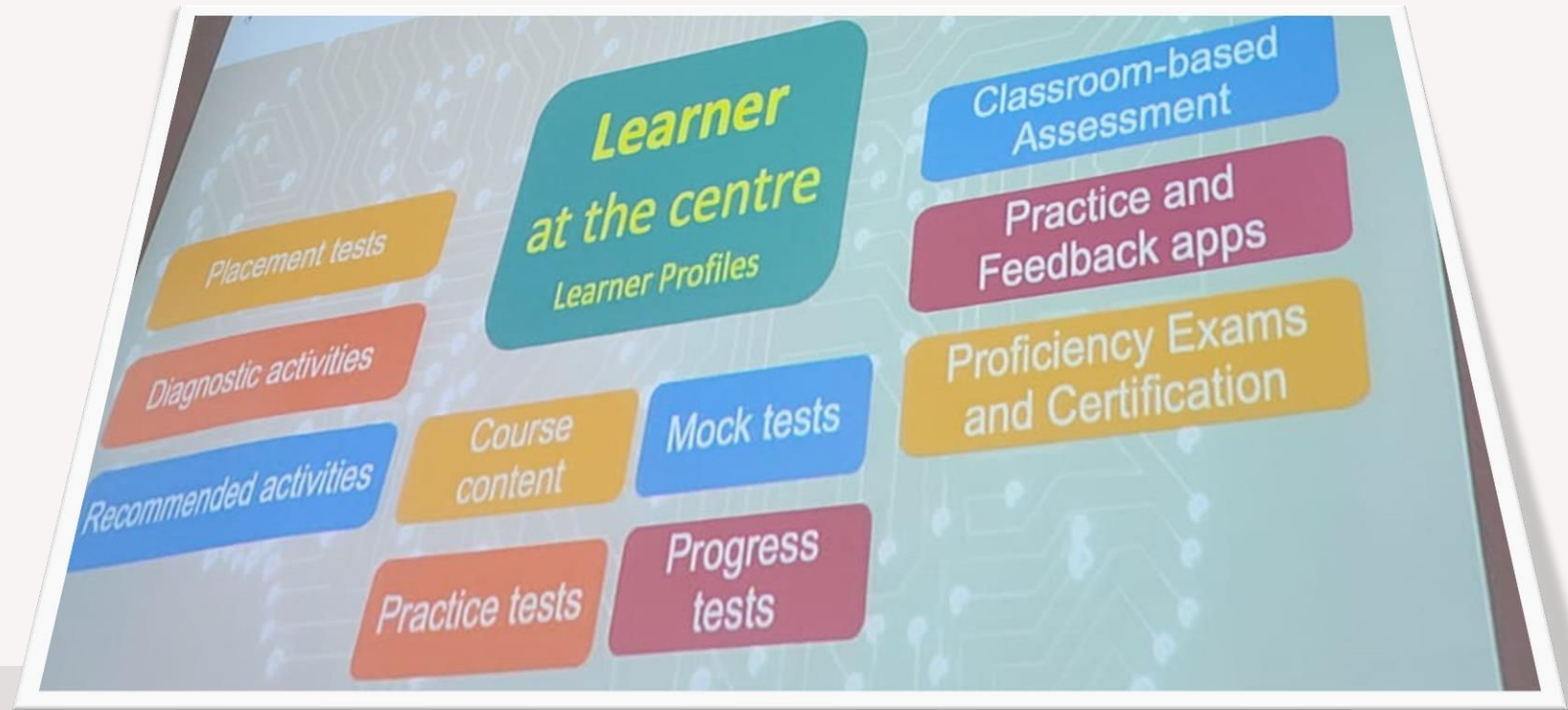
## Retos y desafíos

La IG también plantea algunos retos y desafíos para la evaluación de idiomas, como:

- Precisión:** La precisión de las herramientas de IG es todavía una preocupación, ya que pueden cometer errores que pueden afectar a la validez de las evaluaciones.

- Sesgo:** Las herramientas de IG pueden estar sesgadas en función de los datos en los que se han entrenado, lo que puede afectar a la equidad de las evaluaciones.

# INTEGRACIÓN DE LA EVALUACIÓN Y EL APRENDIZAJE



# IMPACTO

Esta estrategia DE MONITOREO permite al MEP determinar la ubicación lingüística de la población estudiantil de acuerdo al MCER y sugerir recomendaciones curriculares para la mejora de cada modalidad escolar con un instrumento de evaluación personalizado (temas, nivel, y complejidad) con ayuda de la tecnología (IA/IG).



# Gracias!



allenquesada@gmail.com

<https://pelex.ucr.ac.cr>

