



Martin Benjamin & Paula Radetzky

Multilingual Lexicography with a Focus on Less-Resourced Languages:
Data Mining, Expert Input, Crowdsourcing, and Gamification

Workshop on Collaboration and Computing for Under-Resourced Languages in the Linked Open Data Era
Language Resources and Evaluation Conference (LREC 2014)
26 May, 2014 - Reykjavik, Iceland

- **~7000 languages spoken worldwide** (Ethnologue 2014)
- Top 5: Excellent ICT resources
 - English + FIGS (French, Italian, German, Spanish)
- ~30: Significant ICT resources (speech recognition, machine translation, etc)
 - Many official EU languages + Russian, Arabic, Chinese, Japanese...
 - ~~Hindi, Bengali, Punjabi, Indonesian, Swahili...~~ (100 million+ speakers)
- **~80: More than 10,000,000 speakers**
- <100: Basic ICT resources (corpora, machine-readable dictionaries, POS taggers, morphological analyzers, parsers, etc)
- **~300: More than 1,000,000 first language speakers**
- **~1500: More than 100,000 first language speakers**
- **~2000: Threatened – 10,000 to 100,000 speakers**
- ~2000: Some existing data (wordlists, Bibles, recordings)
- **~3000: Endangered**
 - Fewer than 3000 speakers
 - Not being passed to next generation – extinct within 100 years
 - Embedded human cultural patrimony

Multilingual Lexicography with a Focus on **Less-Resourced Languages**:



Multilingual Lexicography with a Focus on Less-Resourced Languages: Data Mining, Expert Input, Crowdsourcing, and Gamification

Workshop on Collaboration and Computing for Under-Resourced Languages in the Linked Open Data Era
Language Resources and Evaluation Conference (LREC 2014)

26 May, 2014 - Reykjavik, Iceland



kamusi is Swahili for *dictionary*



Comprehensive linguistic data for every language



Comprehensive linguistic data for every language

- For people
- For machines



In service since 1994 (originally at Yale Council on African Studies)

International NGO since 2009

- Registered non-profit in USA and Switzerland (and soon Tanzania)

Academic Home since 2013:

EPFL - Swiss Federal Institute of Technology in Lausanne

LSIR - Distributed Systems Information Laboratory



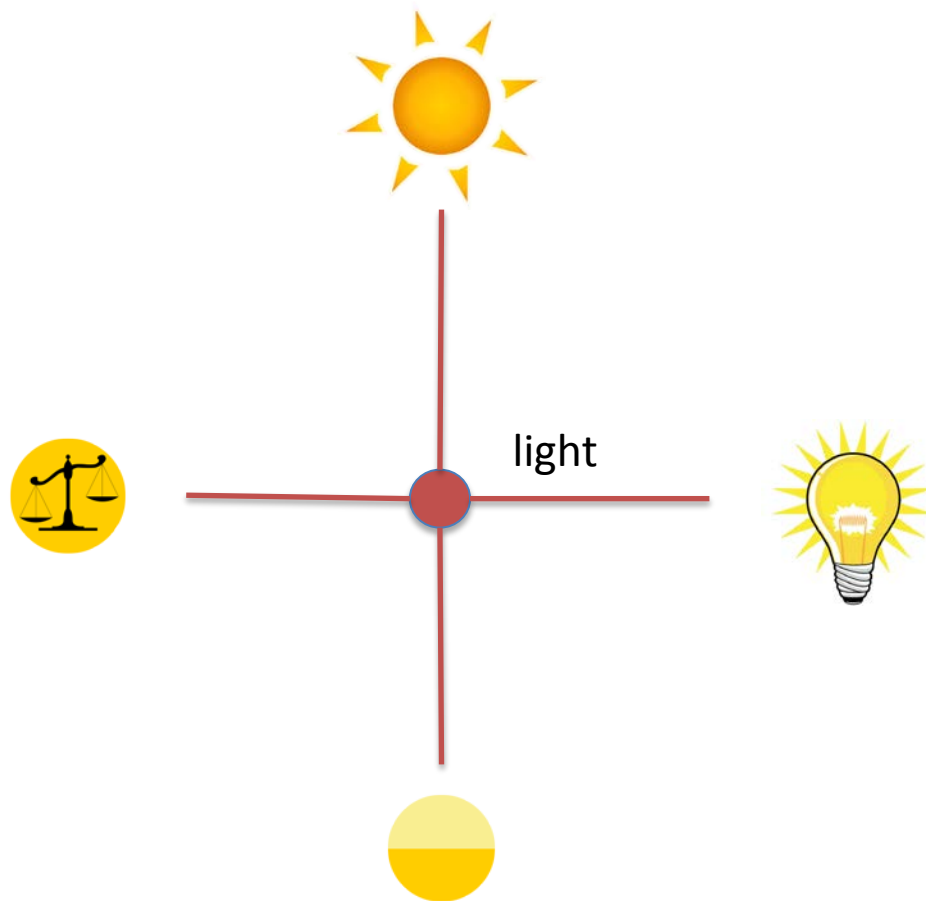
1. Project objectives
2. Acquiring data
3. Crowdsourcing innovations



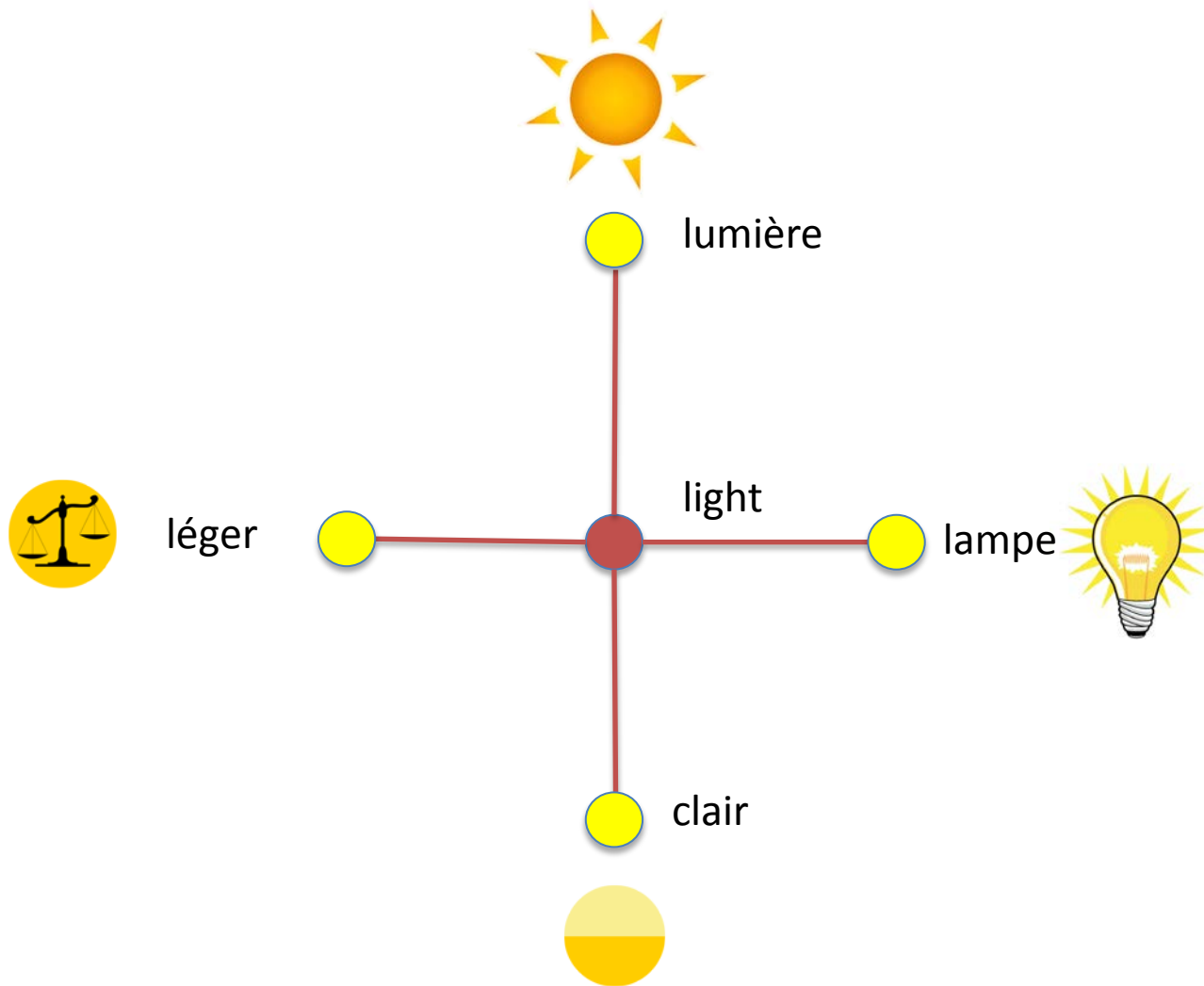
- 1. Project objectives**
2. Acquiring data
3. Crowdsourcing innovations



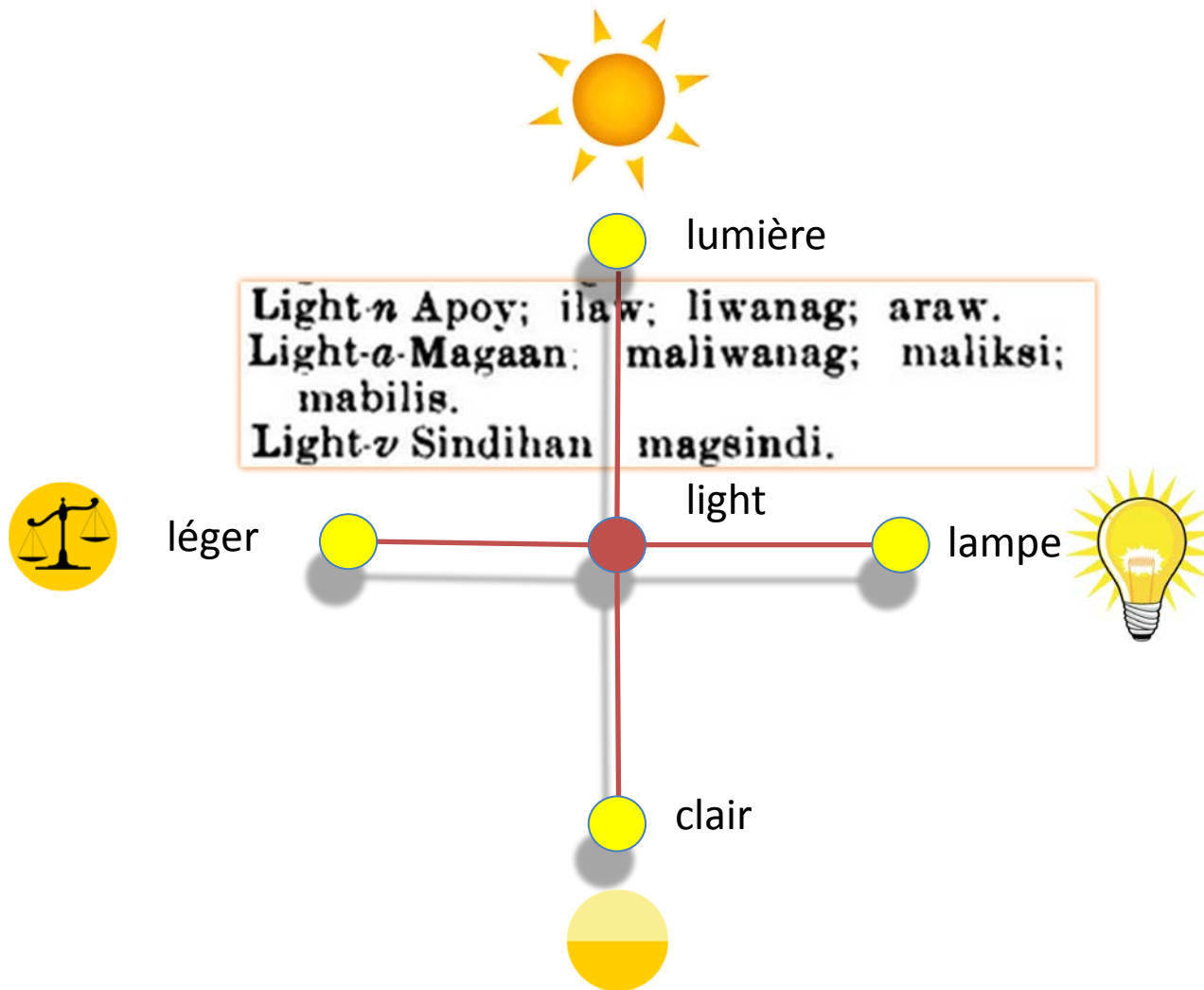
light



why multilingual dictionaries were impossible



why multilingual dictionaries were impossible



why multilingual dictionaries were impossible



light¹ *adj* 1 (of colour) -siokoza, -sioiva angazi, -a ~ brown kabawia isioiva, hafifu. 2 (of a place) -enye mwanga. ~ **coloured** *adj* -enye rangi isioiva. *n* 1 nuru, mwanga *the* ~ *begins to fail* mwanga unaanza kufifia *day* ~ mchana. **in a good/bad** ~ (of picture etc) -a kuonekana vizuri/vibaya; (*fig*) eleweka vizuri/vibaya. **see the** ~ (*liter or rhet*) zaliwa; baini; tangazwa; tambua; -okoka. **be/stand in one's** ~ kinga nuru; (*fig*) zuia mtu anikio/maendeleo kwake. **stand in one's own** ~ zuia kazi yako isionekane; fanya kinyume na matakwa yako. ~ **year** *n* (*astron*) kipimo cha umbali kati ya nyota. 2 taa. ~ **s out** muda wa kuzima taa. **the northern/southern** ~ *n* miali ya mwanga katika ncha za kaskazini na kusini. 3 mwako wa moto; kiberiti *strike a* ~ washa moto; washa kiberiti. 4 uchangamfu (usoni mwa mtu). **the** ~ **of somebody's countenance** (*biblical*) kupendezwa kwake. 5



léger



lampe



why multilingual dictionaries were impossible



in f t ★

ગુજરાતીલેક્સિકોન.કોમ

Ratilal Chandara's Gujarati Language Resources

Gujaratilexicon » Dictionary » English To Gujarati

Dictionary | Opposites | Thesaurus | Idioms | Proverbs | Phrases

Dictionary

English to Gujarati | Gujarati to Gujarati | Gujarati to English | Hindi to Gujarati

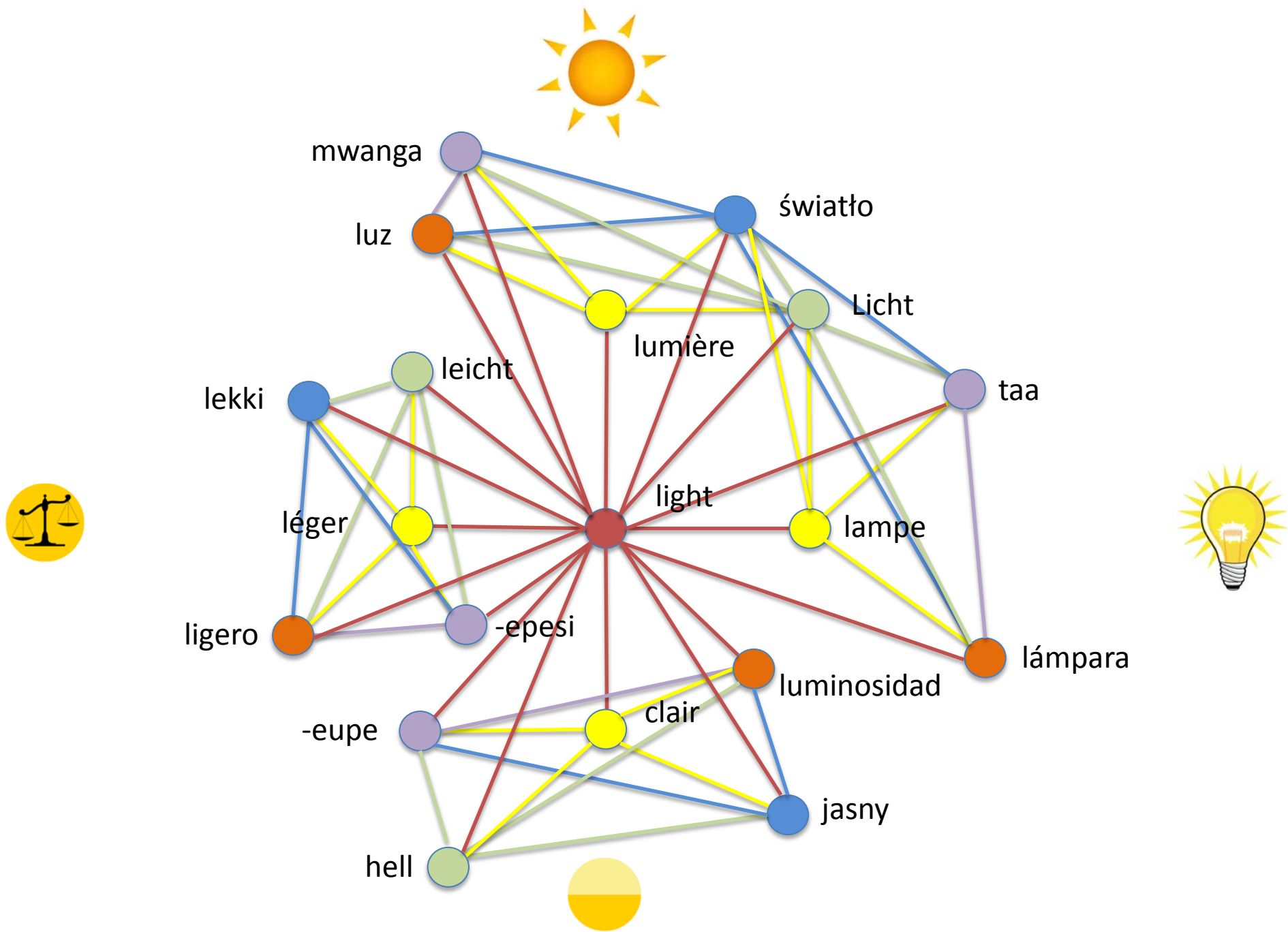
light

Word List

light

No.	Type	Pronunciation	Meaning
1	નિલો	લાઇટ	તેજ, પ્રકાશ, અજ્ઞવાણું, જ્યોતિ, દીપ, દીવો, અવસ્થાવર નિર્મમક દીવો, જેના વડે વસ્તુ દૃશ્યમાન થાય છે તે સાધન-પ્રકાર, દેવતા સમગાવવાની દીવાસળી કે કાકરો, તેજસ્વિતા, આંખનું તેજ, પ્રકાશનું ઘટકોઈ ઉદ્ભવસ્થાન સૂર્ય, મીઠાખતી, ઈ. જેવું, ચમક, કચાકની ઊંજળી ખાજુ, હાકિડોલ, ફાટિ, ઓંધ, સાન, કોઈ ખાખતને સ્પષ્ટ કરનારું સાન-માસિની ઈ., અંધારું-નારિ (રંગ અંગે) ફીકું, પાંખું, ઝાંખું, (આવ ઈ.) સમગાવવું, પ્રકાશવું, પેટવડાવવું, અગતું, દીવાથી રસનો બતાવવો, -ને પ્રકાશ આપવો, પ્રસન્ન થવું કે કરવું, ઉદ્દાસિત થવું કે કરવું

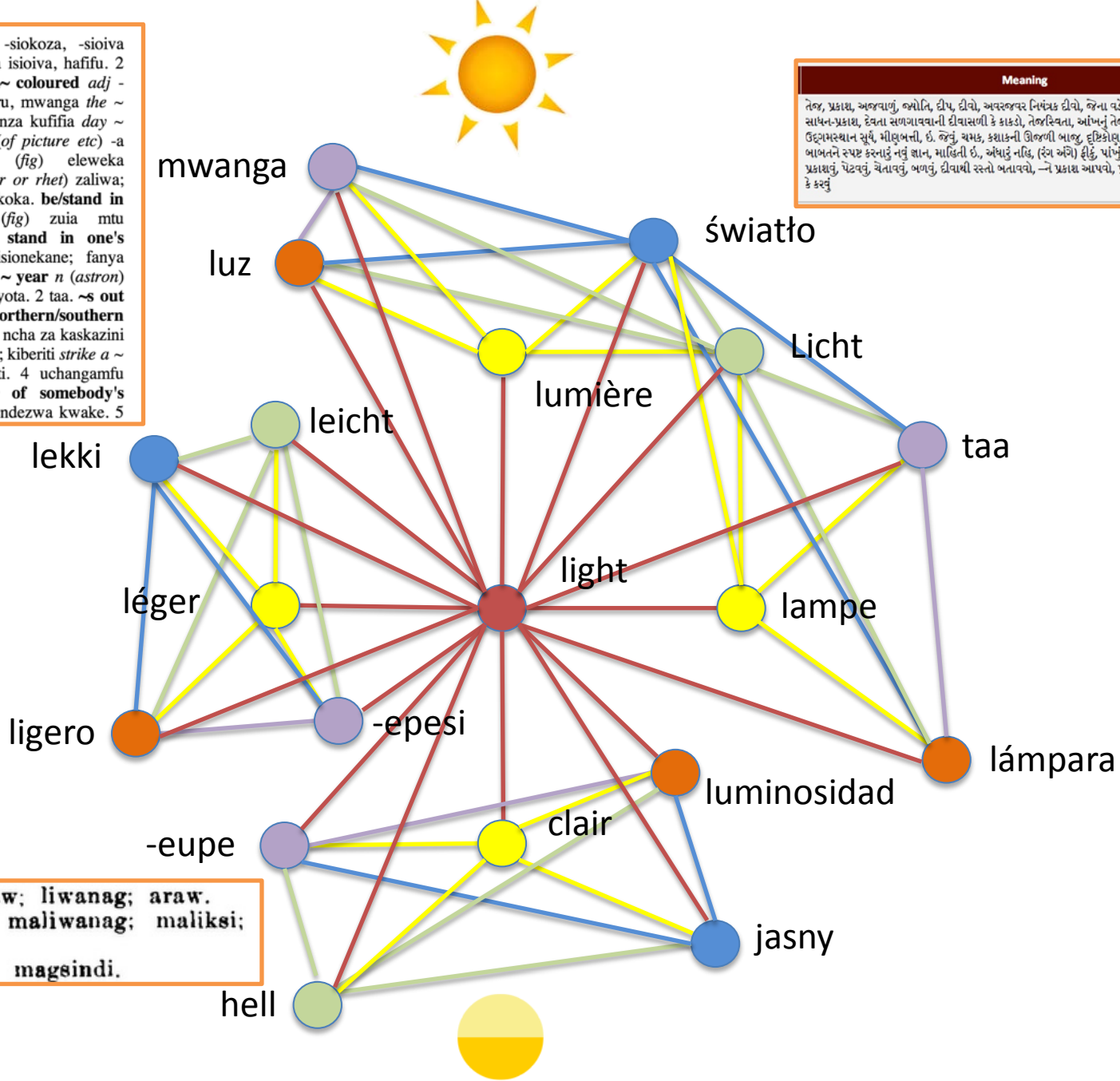
why multilingual dictionaries were impossible



why multilingual dictionaries were impossible

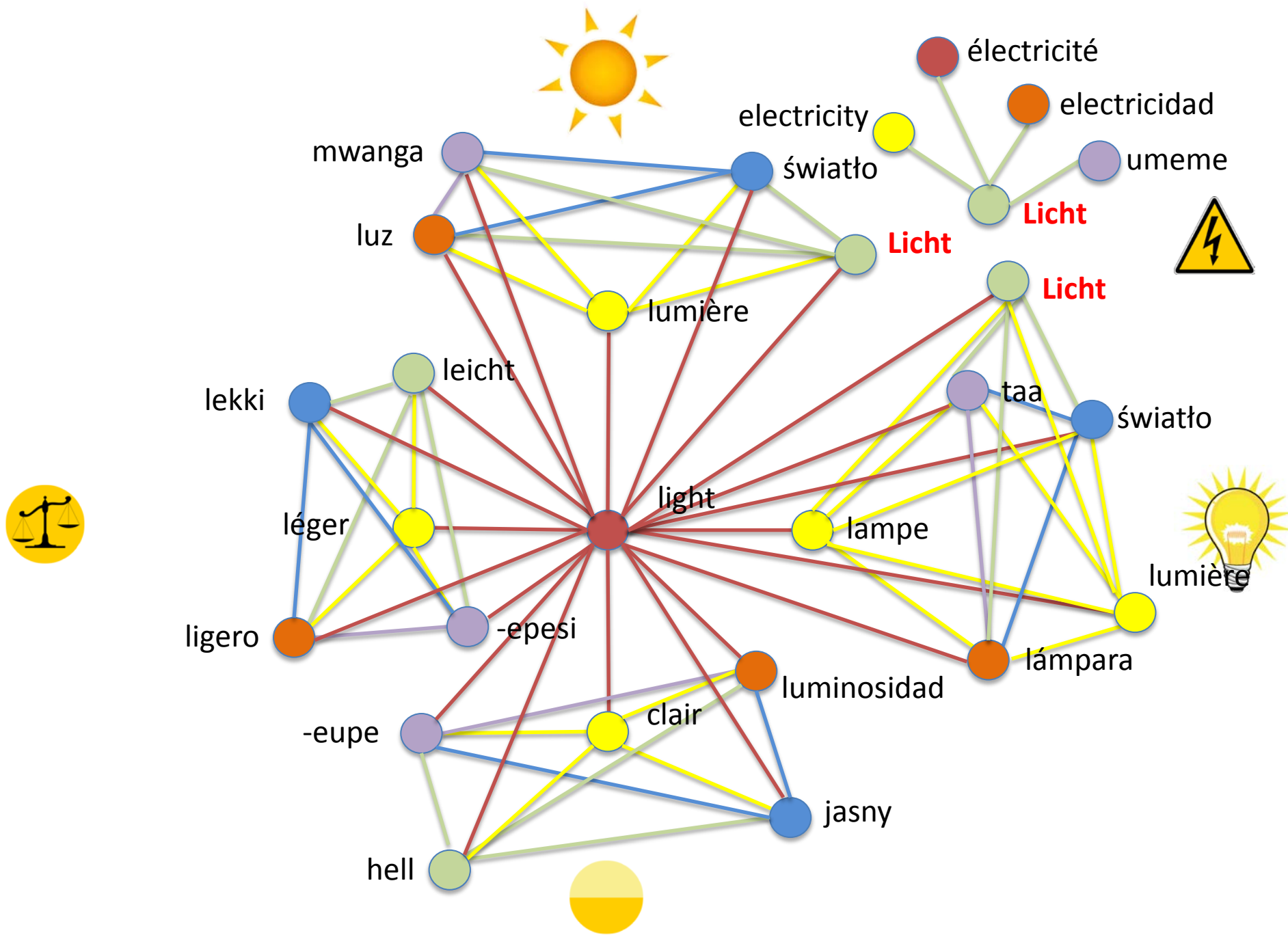
light¹ *adj* 1 (of colour) -siokoza, -sioiva angazi, -a ~ brown kahawia isioiva, hafifu. 2 (of a place) -enye mwanga. ~ **coloured** *adj* -enye rangi isioiva. *n* 1 nuru, mwanga *the* ~ *begins to fail* mwanga unaanza kufifia *day* ~ mchana. **in a good/bad** ~ (of picture etc) -a kuonekana vizuri/vibaya; (*fig*) eleweka vizuri/vibaya. **see the** ~ (*liter or rhet*) zaliwa; baini; tangazwa; tambua; -okoka. **be/stand in one's** ~ kinga nuru; (*fig*) zuia mtu mafanikio/maendeleo yake. **stand in one's own** ~ zuia kazi yako isionekane; fanya kinyume na matakwa yako. ~ **year** *n* (*astron*) kipimo cha umbali kati ya nyota. 2 taa. ~ **s out** muda wa kuzima taa. **the northern/southern** ~ *n* miali ya mwanga katika ncha za kaskazini na kusini. 3 mwako wa moto; kiberiti *strike a* ~ washa moto; washa kiberiti. 4 uchangamfu (usoni mwa mtu). **the** ~ of somebody's countenance (*biblical*) kupendezwa kwake. 5

Meaning
 तेज, प्रकाश, अलखवाणुं, ज्योति, दीप, दीवो, अवलोकनं निबन्धक दीवो, ज्योति वदो वस्तु ह्ययम्मानं यद्य छे ते साधन-प्रकाश, देवता सणजावचानी दीवासणी के काकरो, तेजस्विता, आभिननुं तेज, प्रकाशनुं वस्कोठं उद्गमस्थानं सूर्य, मीलुभनी, छं, ज्युं, यमक, कशाकनी जिल्ली भाजू, दृष्टिकोष, दृष्टि, भोध, ज्ञान, कोठं आभनने रूपक करानुं ननुं ज्ञान, मासिनी छं, अंधानुं नहि, (रंज अंगि) हीकुं, पांभुं, जंभुं, (आज छं.) सणजावचुं, प्रकाशनुं, पेटवचुं, येतावचुं, अणनुं, दीवामी रस्तो अतावचो, -ने प्रकाश आपचो, प्रसस यनुं के करुं, उलसित यनुं के करुं

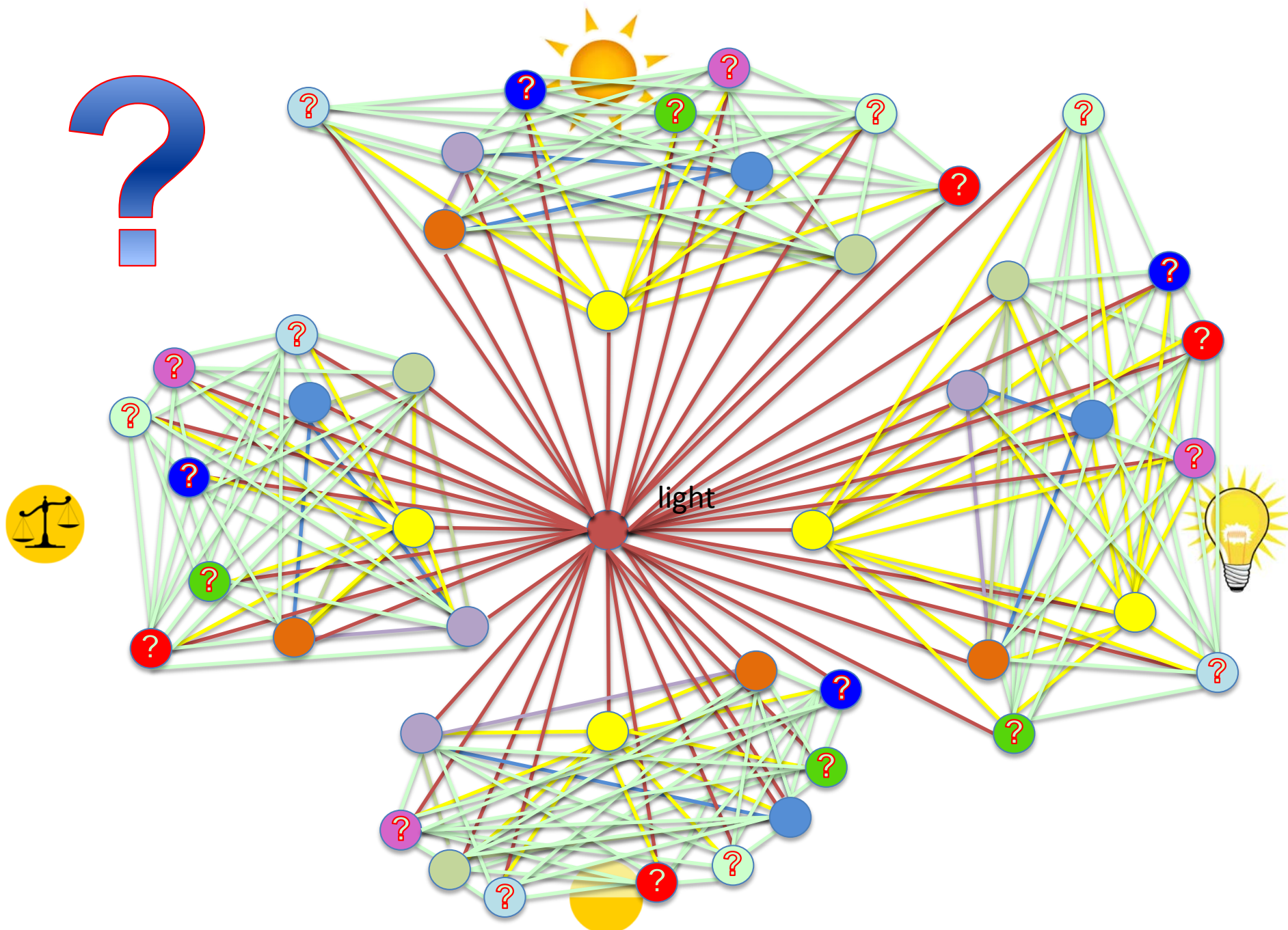


Light *n* Apoy; ilaw; liwanag; araw.
Light *a* Magaan; maliwanag; maliksi; mabilis.
Light *v* Sindihan magsindi.

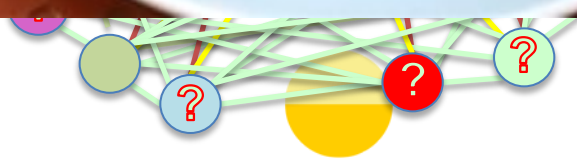
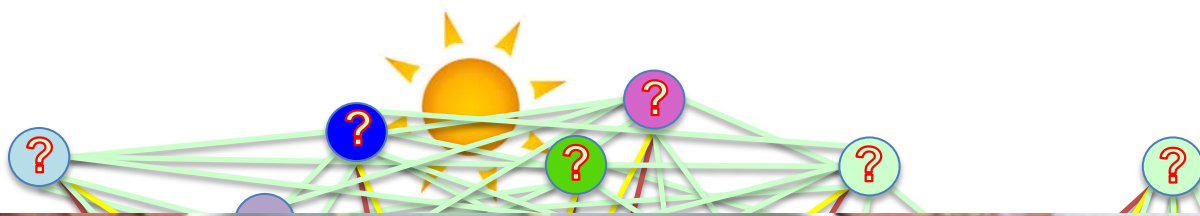
why multilingual dictionaries were impossible



why multilingual dictionaries were impossible



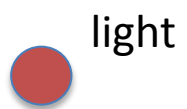
why multilingual dictionaries based on spelling are impossible



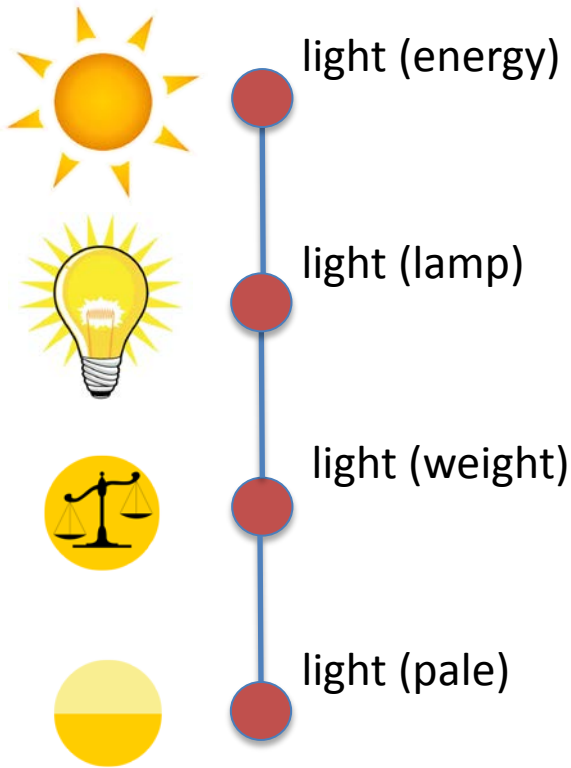
why multilingual dictionaries based on spelling are impossible



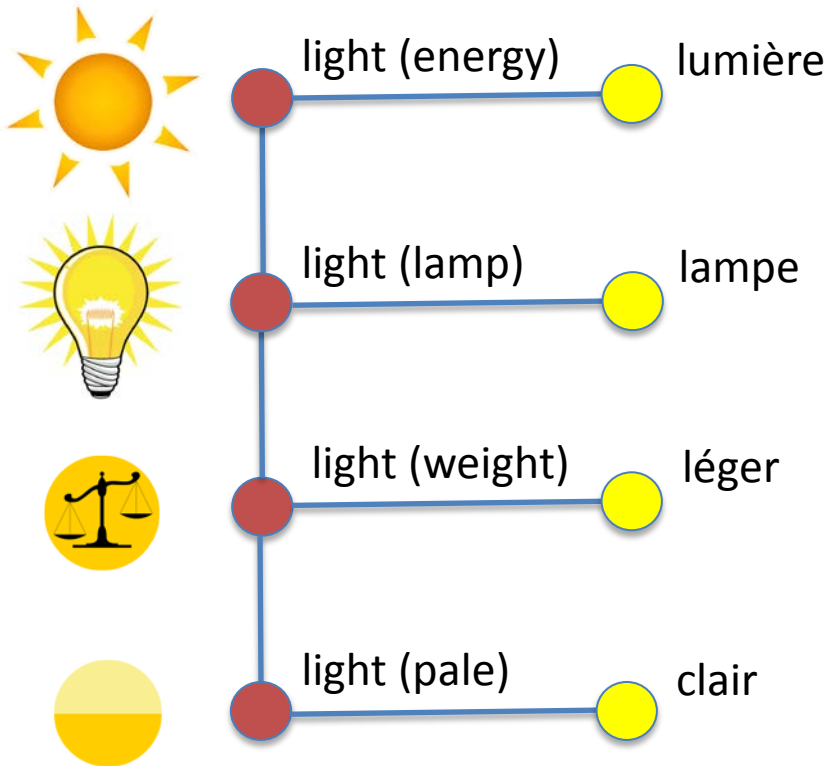
impossible



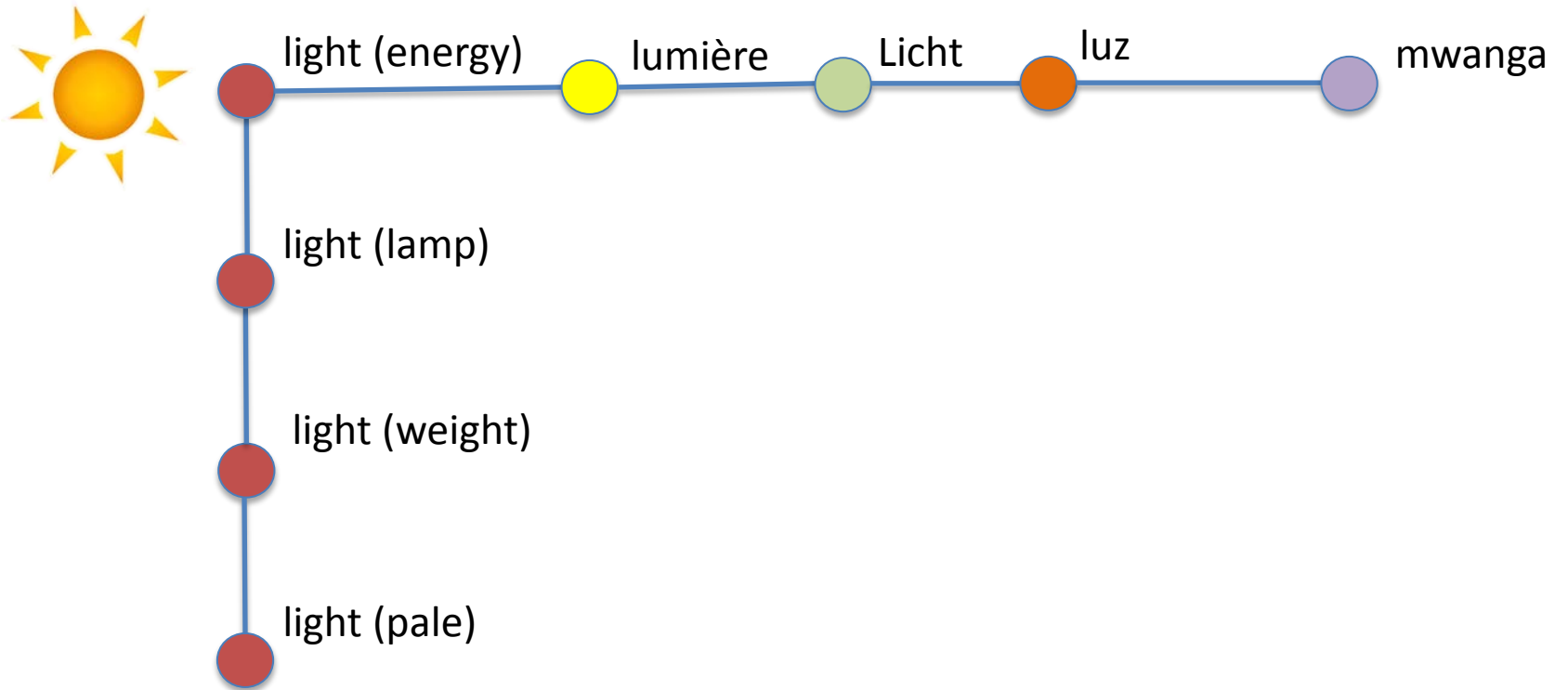
how Kamusi makes a multilingual dictionary **possible**



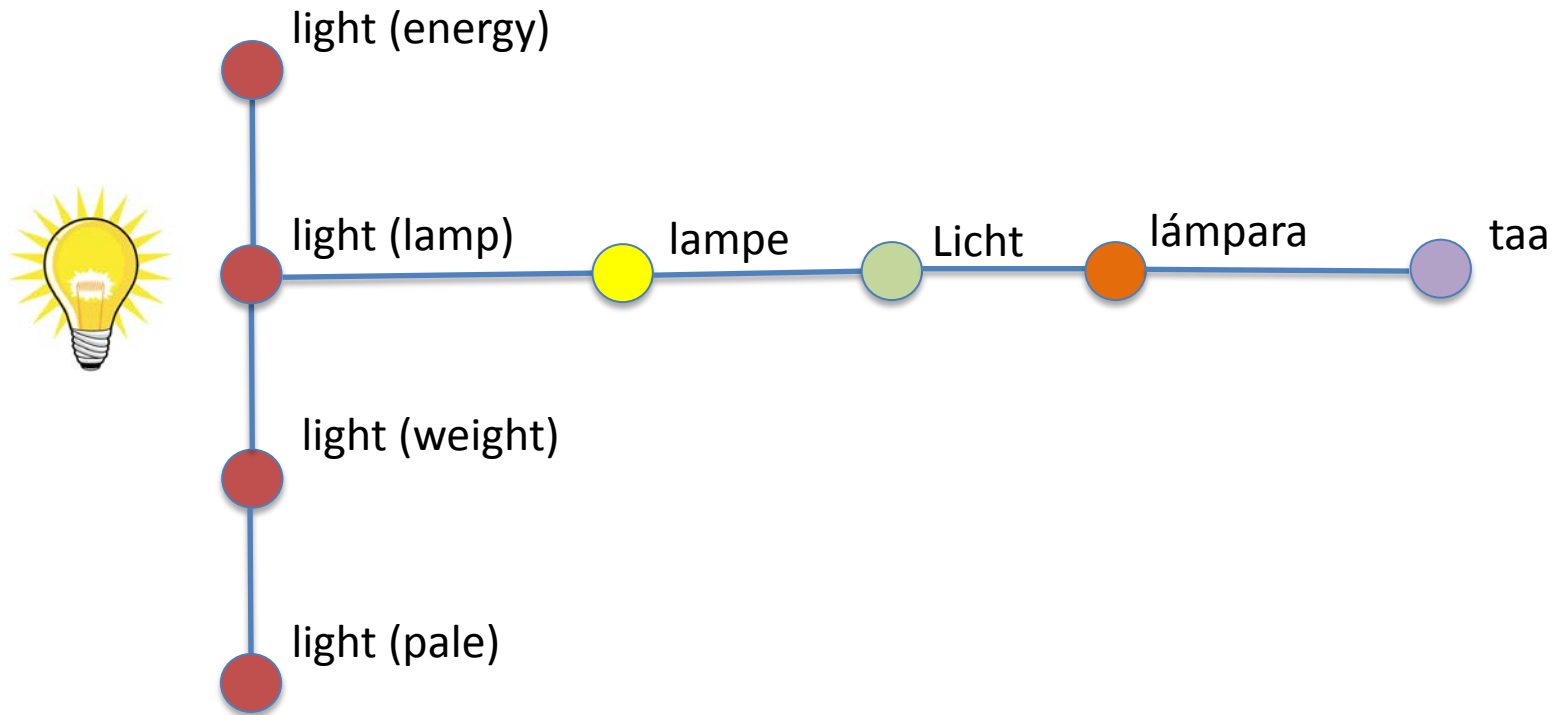
how Kamusi makes a multilingual dictionary possible



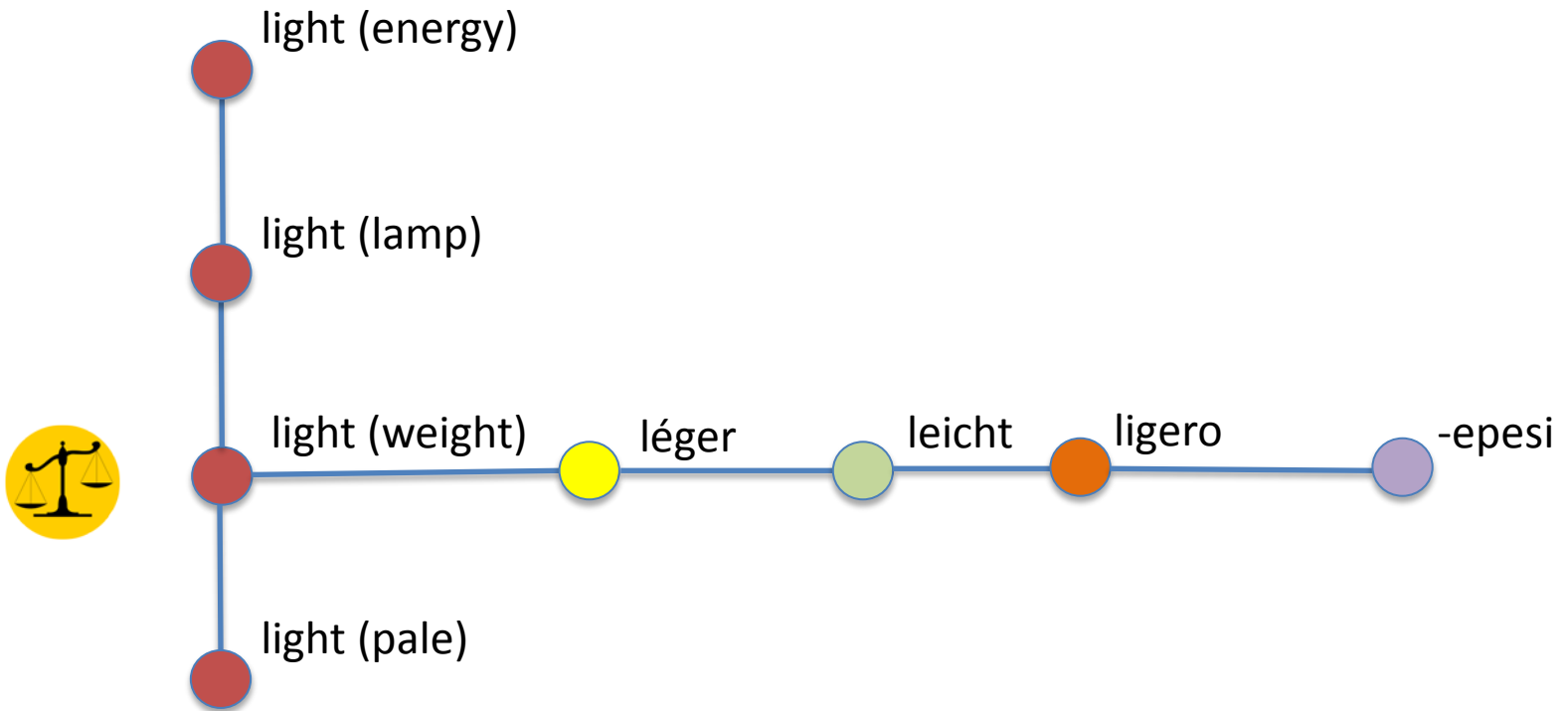
how Kamusi makes a multilingual dictionary possible



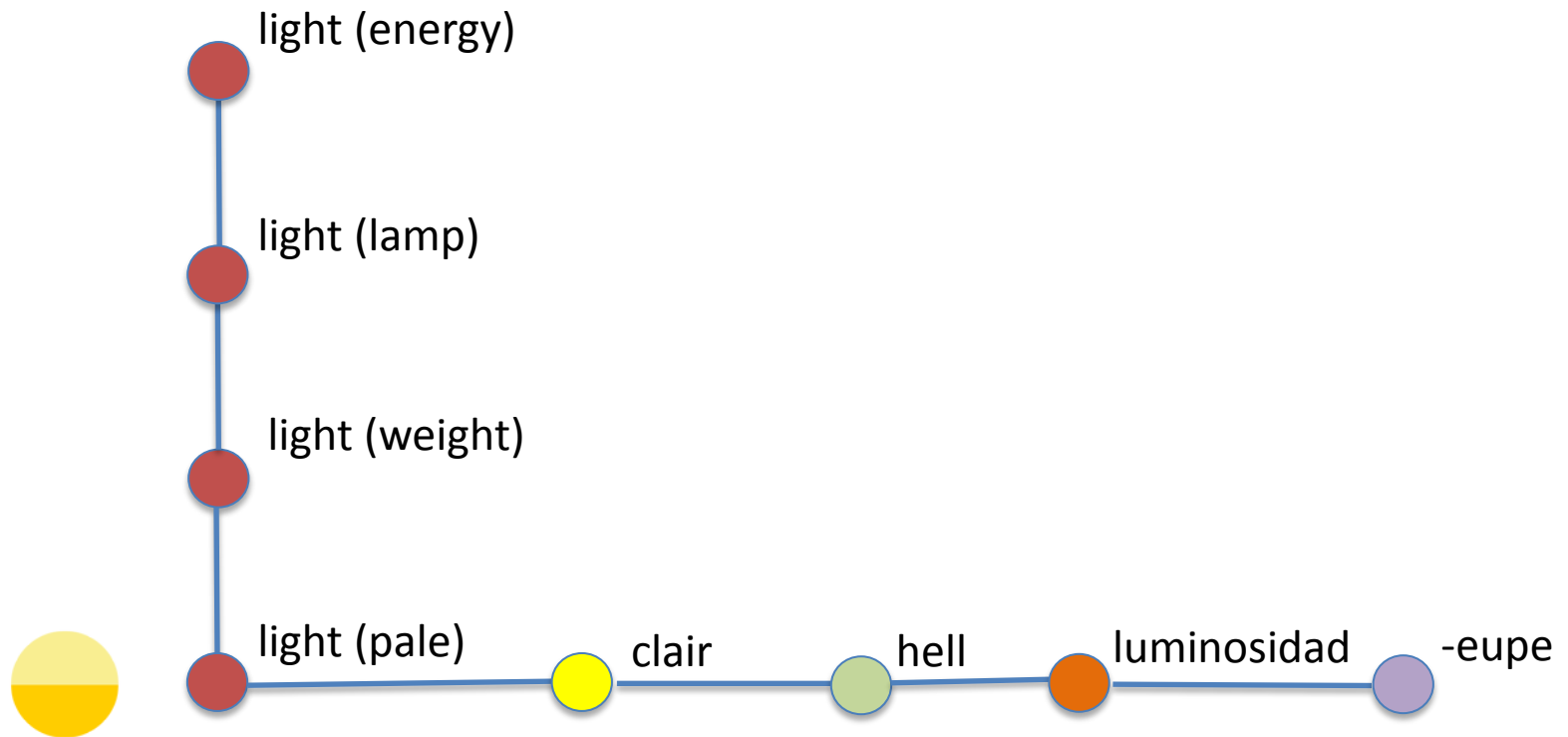
how Kamusi makes a multilingual dictionary possible



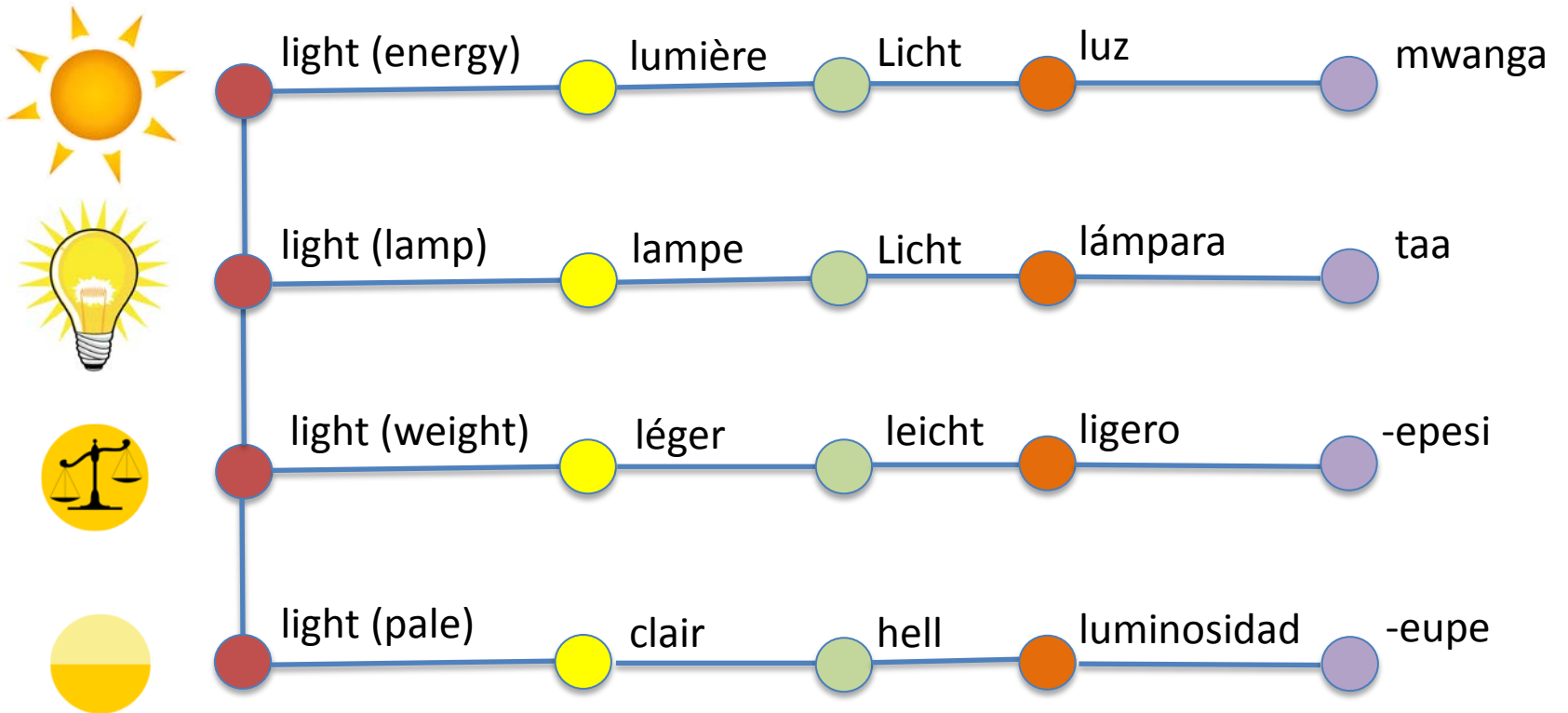
how Kamusi makes a multilingual dictionary possible



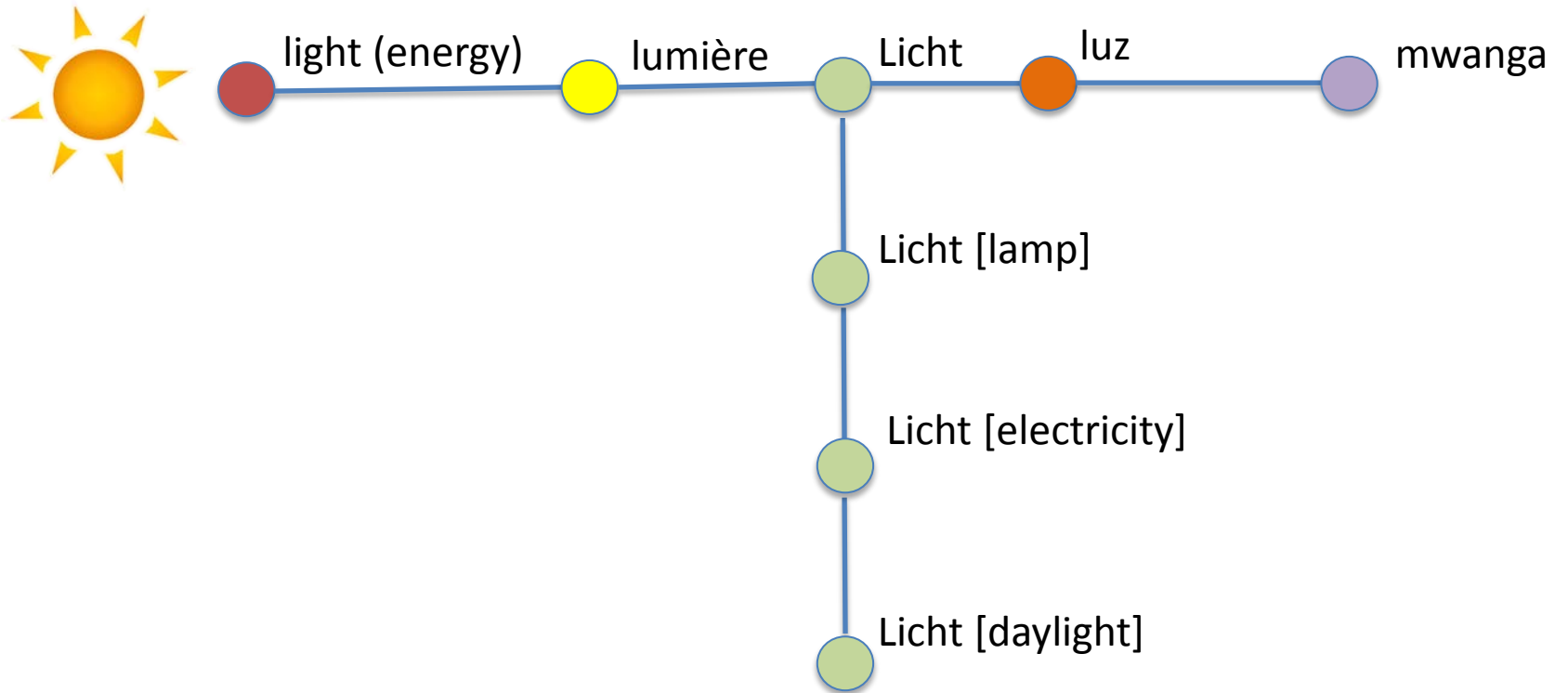
how Kamusi makes a multilingual dictionary possible



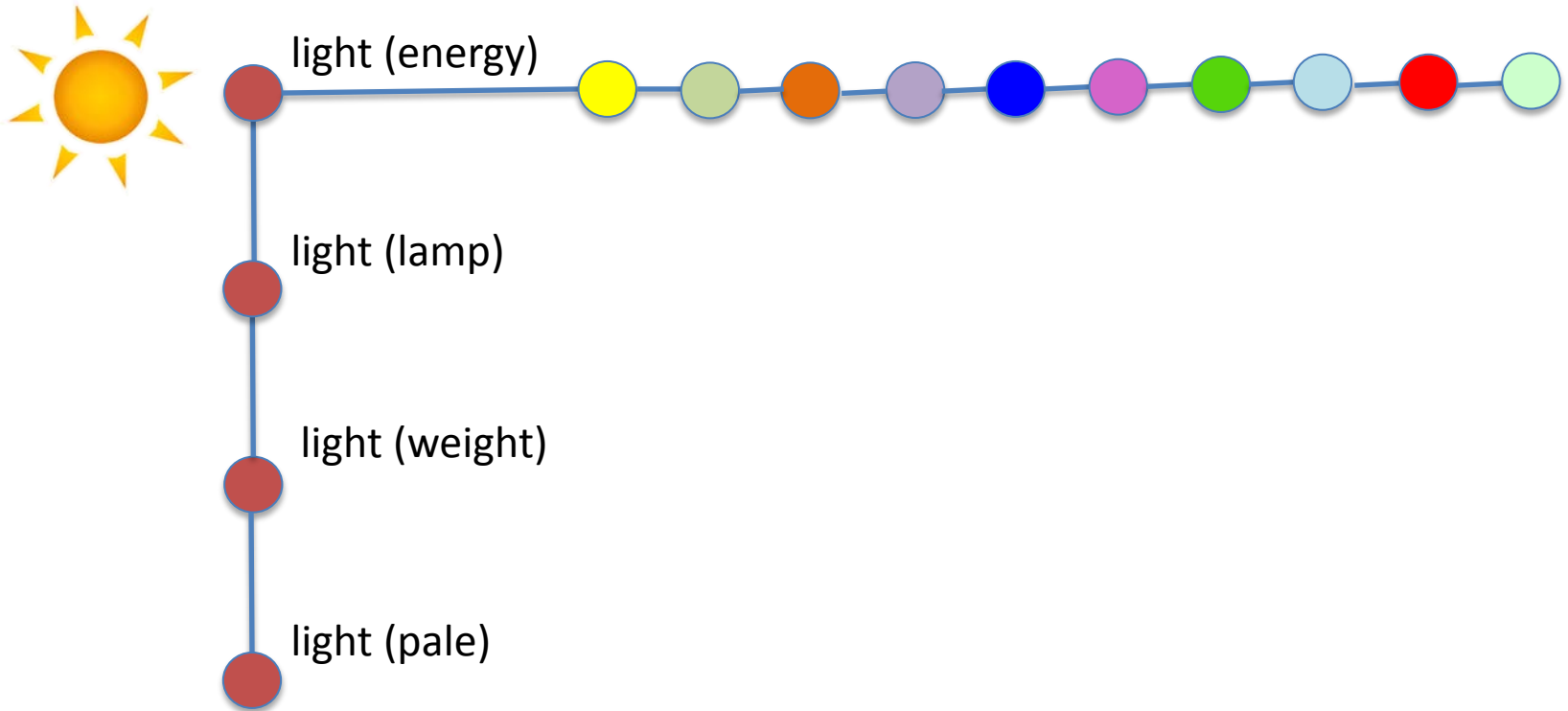
how Kamusi makes a multilingual dictionary possible



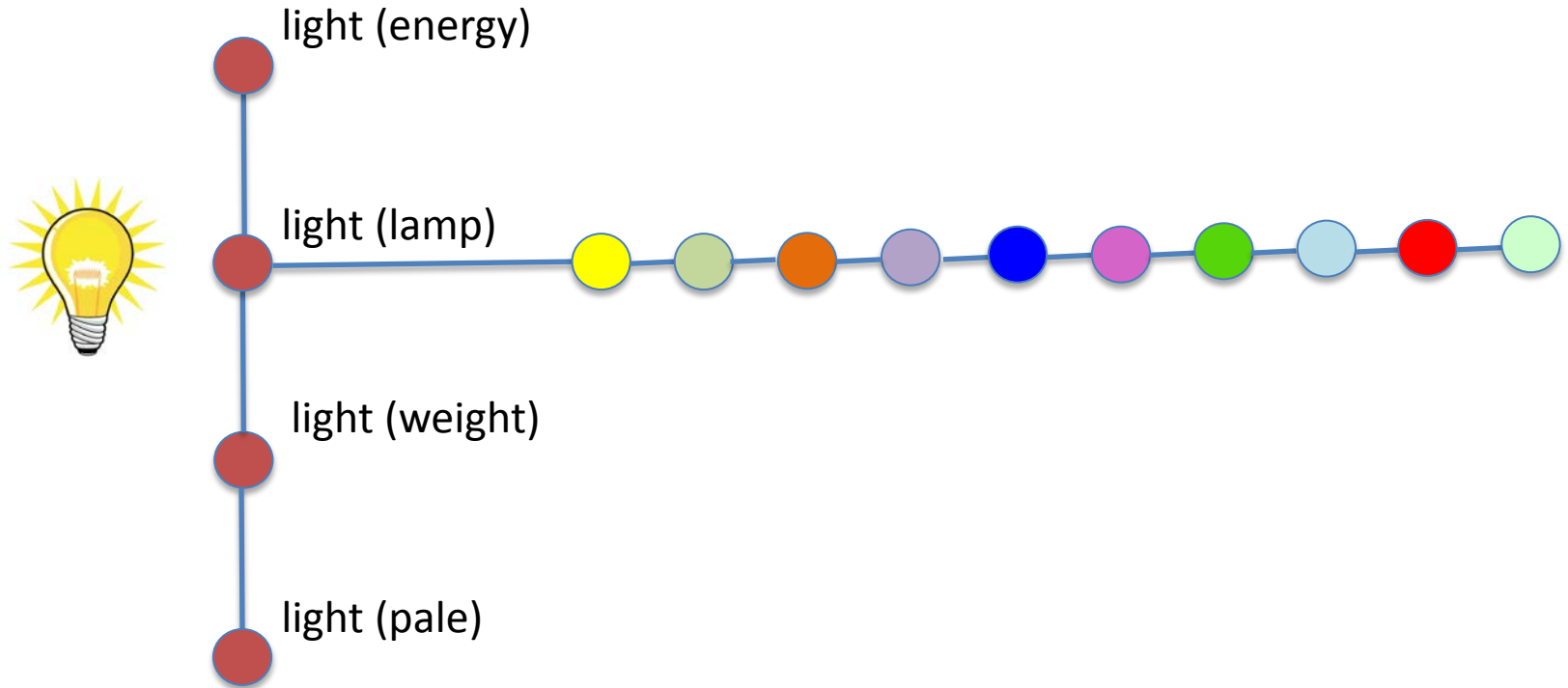
how Kamusi makes a multilingual dictionary possible



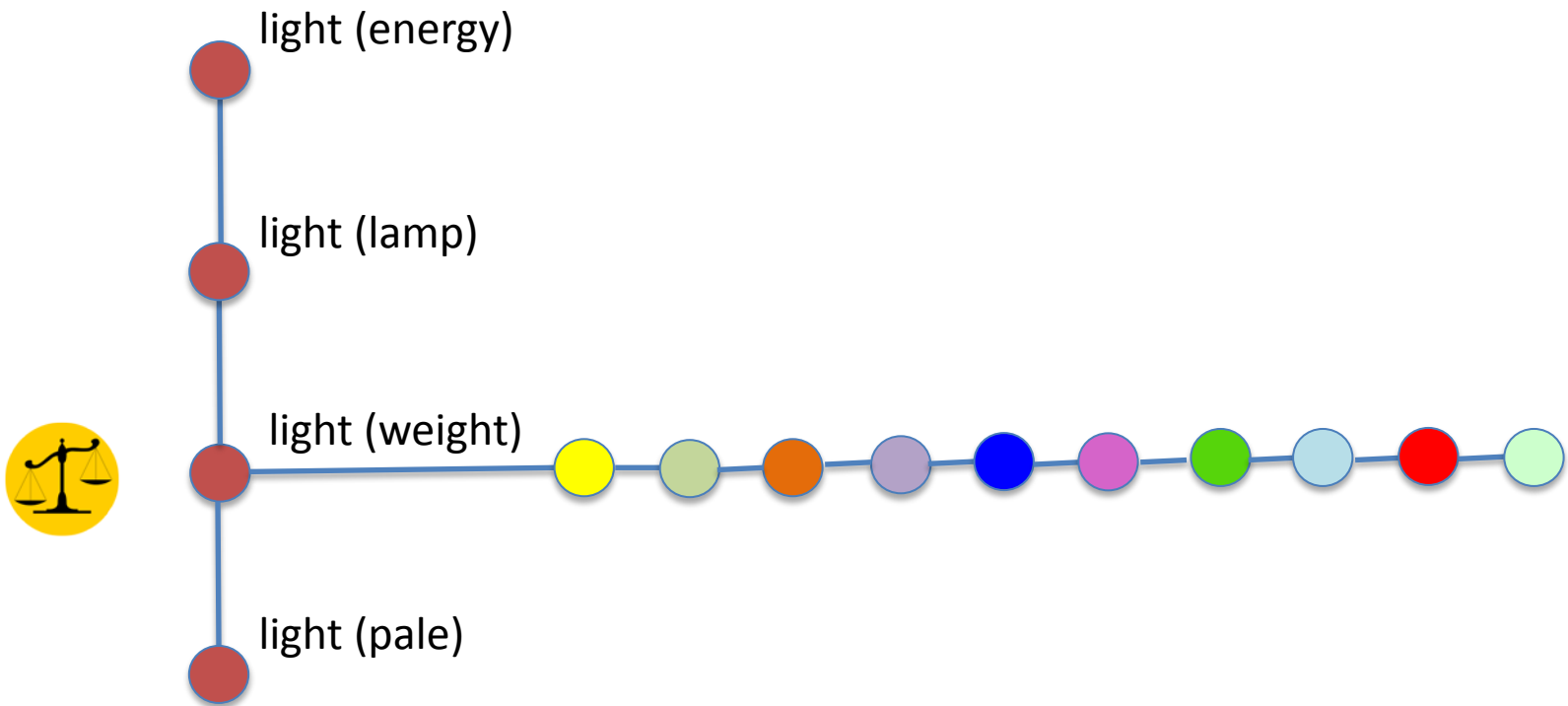
how Kamusi makes a multilingual dictionary possible



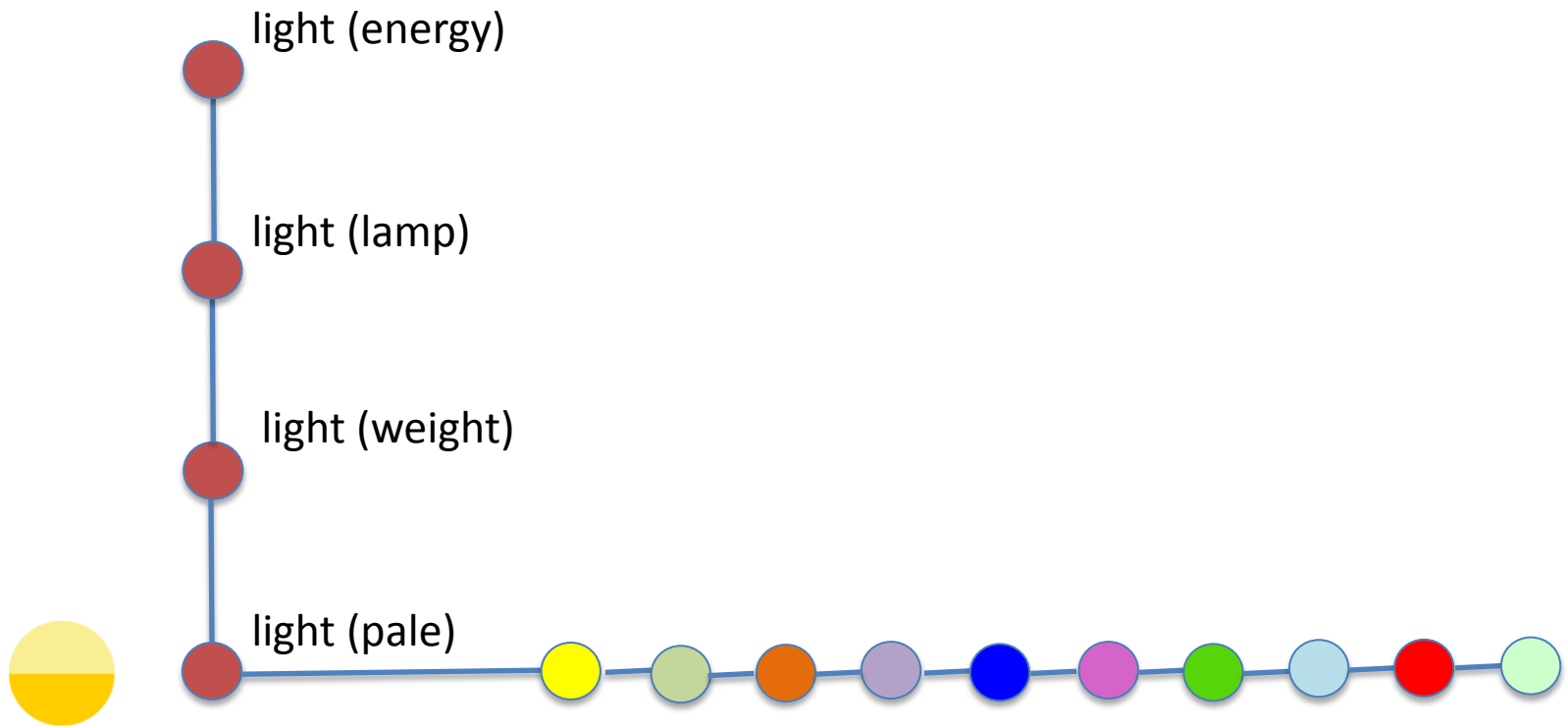
how Kamusi makes a multilingual dictionary possible



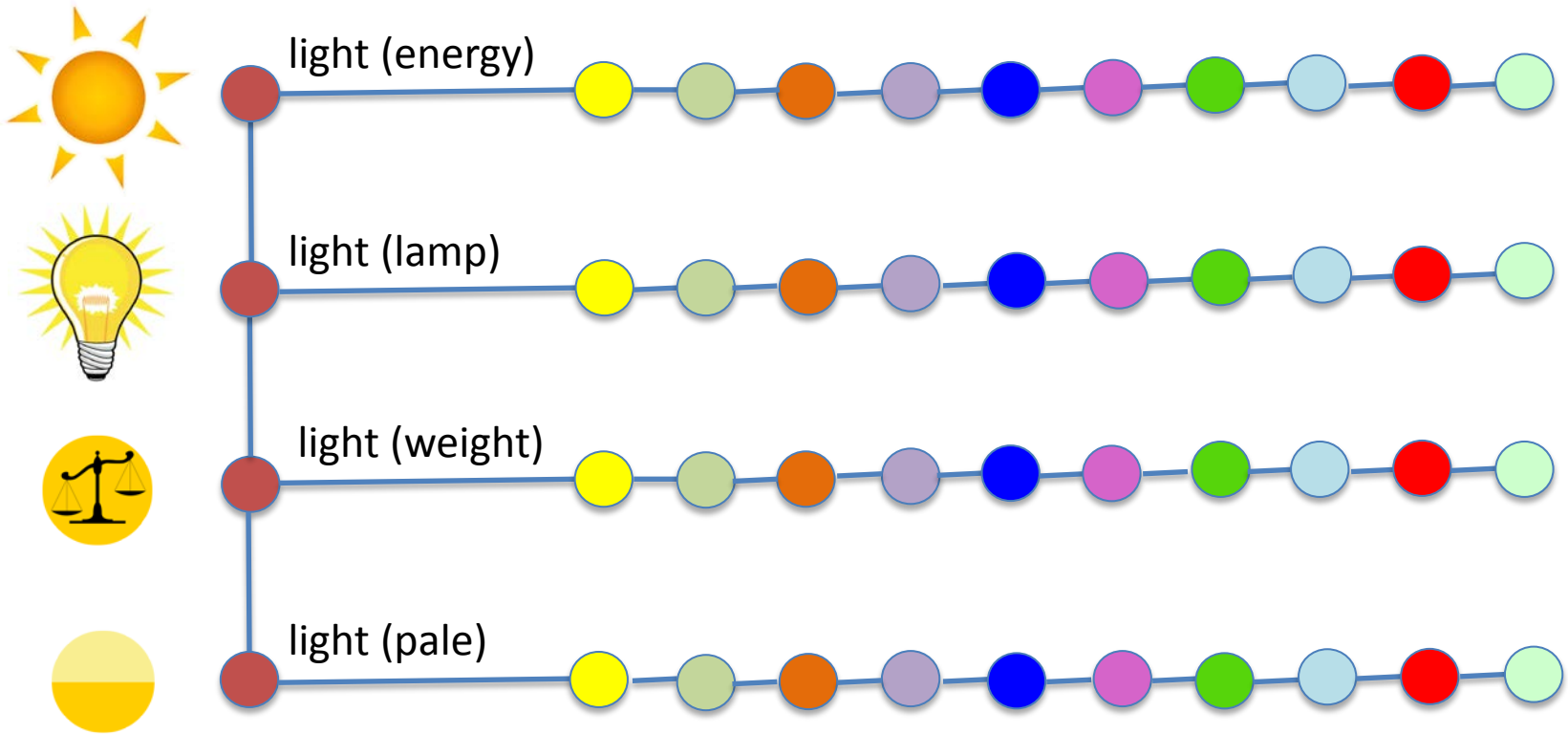
how Kamusi makes a multilingual dictionary possible



how Kamusi makes a multilingual dictionary possible

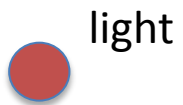
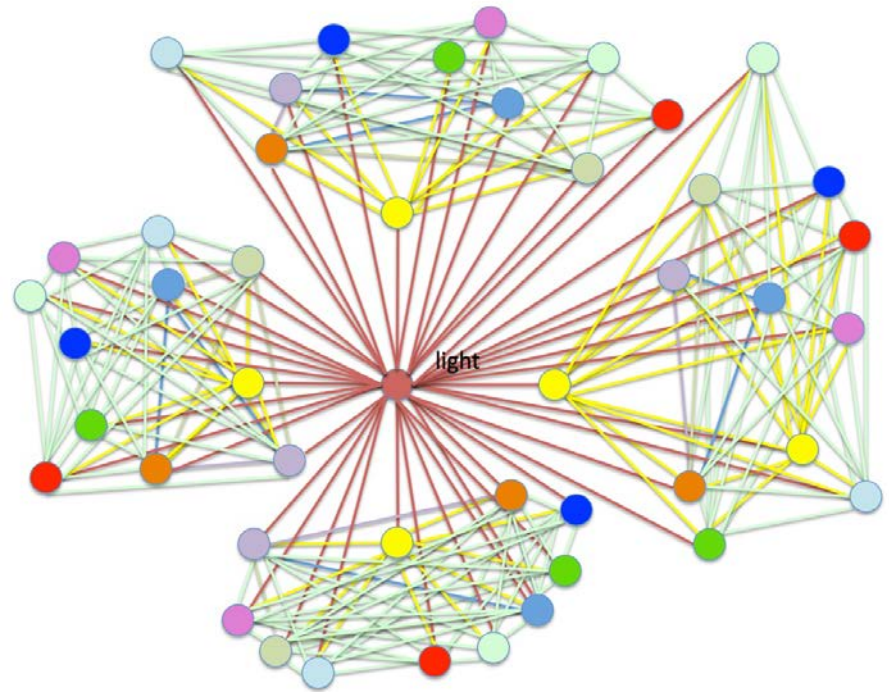


how Kamusi makes a multilingual dictionary possible

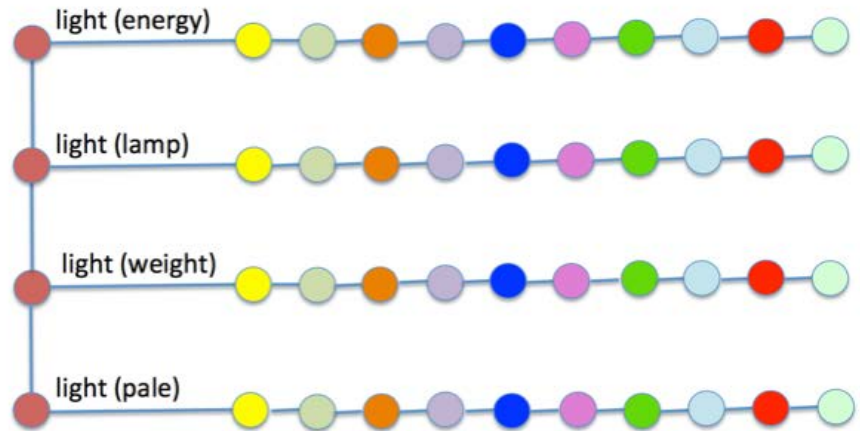


how Kamusi concept linking makes a multilingual dictionary possible

impossible



possible



light

View current

Edit latest

Revisions

Devel

light



View in context



Add an example



Add a translation

English noun

Definition: A device that provides illumination by burning electricity, fuel, or wax.

Terminology: general

Plural

lights

My Languages

- **itara** in Kirundi
- **taa** in Swahili
- **lumière** in French
- **lumină** in Romanian

Other Languages

- **taala** in Luganda

licht

View current

Edit current

Revisions

Devel

licht



View in context



Add an example



Add a translation

Dutch zelfstandig naamwoord - noun (Dutch) het-woord, onzijdig (het)

Definition: Voorwerp dat licht produceert, gevoed door elektriciteit of door verbranding van was of een andere brandstof.

Examples:

We gaan slapen. Doe jij het licht (de lichten) uit?

meervoud onzijdig (de)

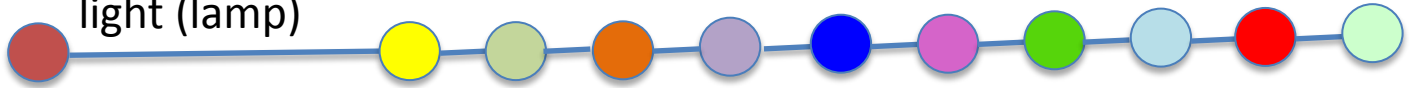
lichten

Other Languages

- **lámpara** in Spanish
- **bugum** in Mampruli
- **-tara** in Kinyarwanda
- **بتي** in Urdu
- **oborabu** in Gusii
- **Licht** in German



light (lamp)



lumină

→ 1° **lumière** (French)

→ 2° **light** (English)

→ 3° **taa** (Swahili)

→ 3° **Licht** (German)

→ 4° **światło** (Polish)





monolingual data

depth: kamusi (unlimited)

Key Information	Word Forms & Origins	Translations & Concepts	Related Intralanguage	Usage Examples
Language	Phonetic IPA	Gloss	Synonyms	Example 1
Lemma	Alt Spelling	Gloss Language	Antonyms	Source
Headword	Tone Spelling	Relation Type	Spawn	Audio
Part of Speech	Alt Script	Science Taxon	Family	Translation
Morphemes	Pronunciations	Term Set	Misspellings	Trans Lang
Definition	Pron Geo-Tag	Tags		Example X
Def Trans	Dialect	Semantic Sets		Source
Def Trans Lang	Word Geo-Tag	Image		Audio
Usage Note	Etym Term	Concept Links		Translation
Cultural Note	Etym Lang			Trans Lang
Special Note	Etym Note			

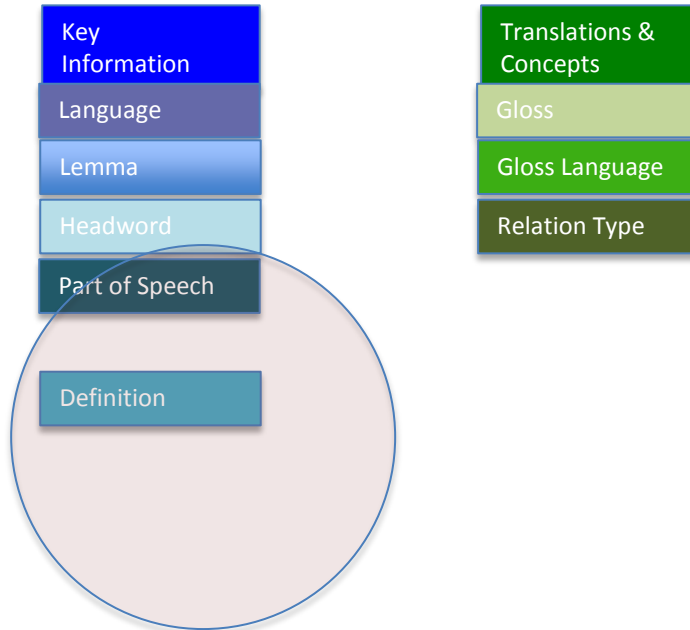
growth: a “living” dictionary

Key Information
Language
Lemma

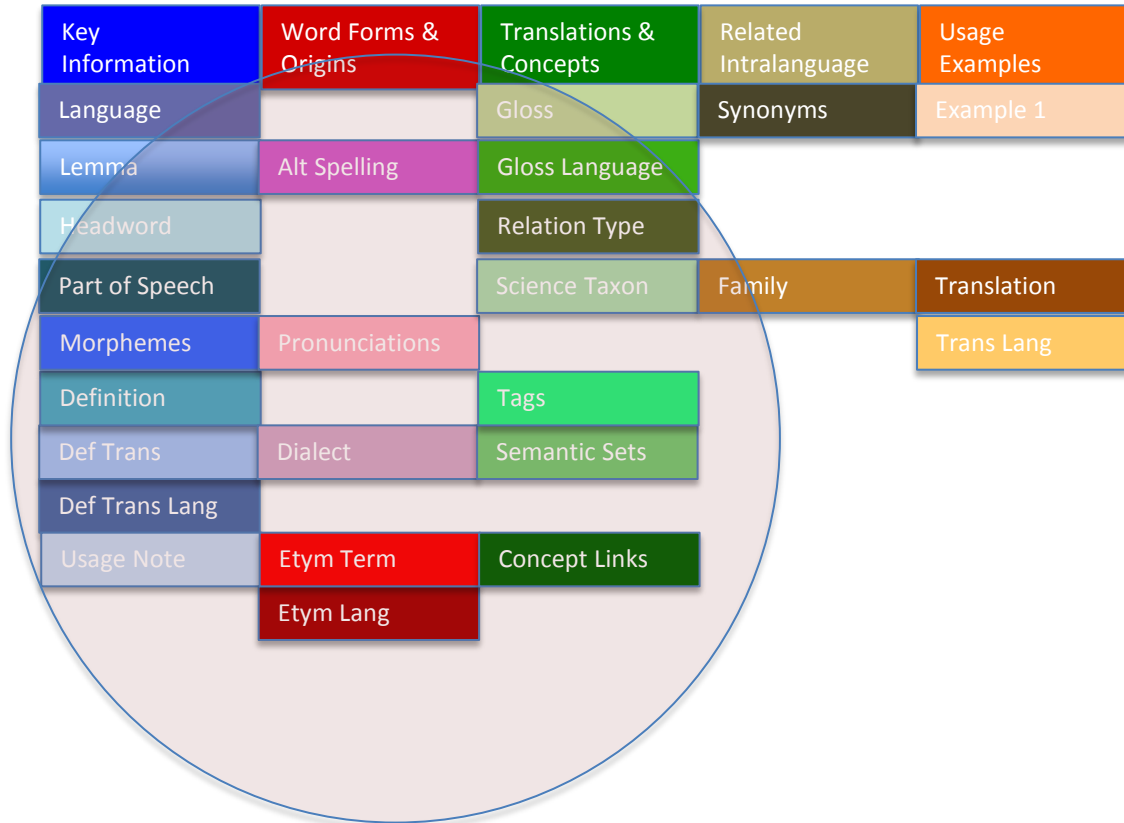
Translations & Concepts
Gloss
Gloss Language



growth: a “living” dictionary



growth: a “living” dictionary



growth: a “living” dictionary

Key Information	Word Forms & Origins	Translations & Concepts	Related Intralanguage	Usage Examples
Language	Phonetic IPA	Gloss	Synonyms	Example 1
Lemma	Alt Spelling	Gloss Language	Antonyms	Source
Headword	Tone Spelling	Relation Type	Spawn	Audio
Part of Speech	Alt Script	Science Taxon	Family	Translation
Morphemes	Pronunciations	Term Set	Misspellings	Trans Lang
Definition	Pron Geo-Tag	Tags		Example X
Def Trans	Dialect	Semantic Sets		Source
Def Trans Lang	Word Geo-Tag	Image		Audio
Usage Note	Etym Term	Concept Links		Translation
Cultural Note	Etym Lang			Trans Lang
Special Note	Etym Note			

Edit Dictionary Term licht

View current

Edit current

Revisions

Devel

Key Information *

Word Forms and Origins

Translations and Concepts

Related terms within the same language

Examples *

Concept

Term (lemma) *

licht

Term Language *

Dutch

Headword *

(?) ^o licht

Part of Speech *

zelfstandig naamwoord - noun

Part Attribute

het-woord, onzijdig (het)

Morpheme Information

Collected Morphemes

meervoud mannelijk (de)

+

enkelvoud mannelijk de-woord (de-m)

+

enkelvoud onzijdig het-woord (het)

+

enkelvoud vrouwelijk de-woord (de-v)

+

meervoud vrouwelijk (de)

+

meervoud onzijdig (de)

+

lichten [nid:263256]

Definition

Definition of the term in the language to which it belongs

(?) ^o

Voorwerp dat licht produceert, gevoed door elektriciteit of door verbranding van was of een andere brandstof.

Definition Source URL

Definition Translations

and more...

- Morphemes and inflections
- Attributes
- Behaviors
- Grouping and ranking
- Agglutination parsing
- Grammatical bridges
- Family relations
- Etymology and temporality
- Geography
- And more...

hooks



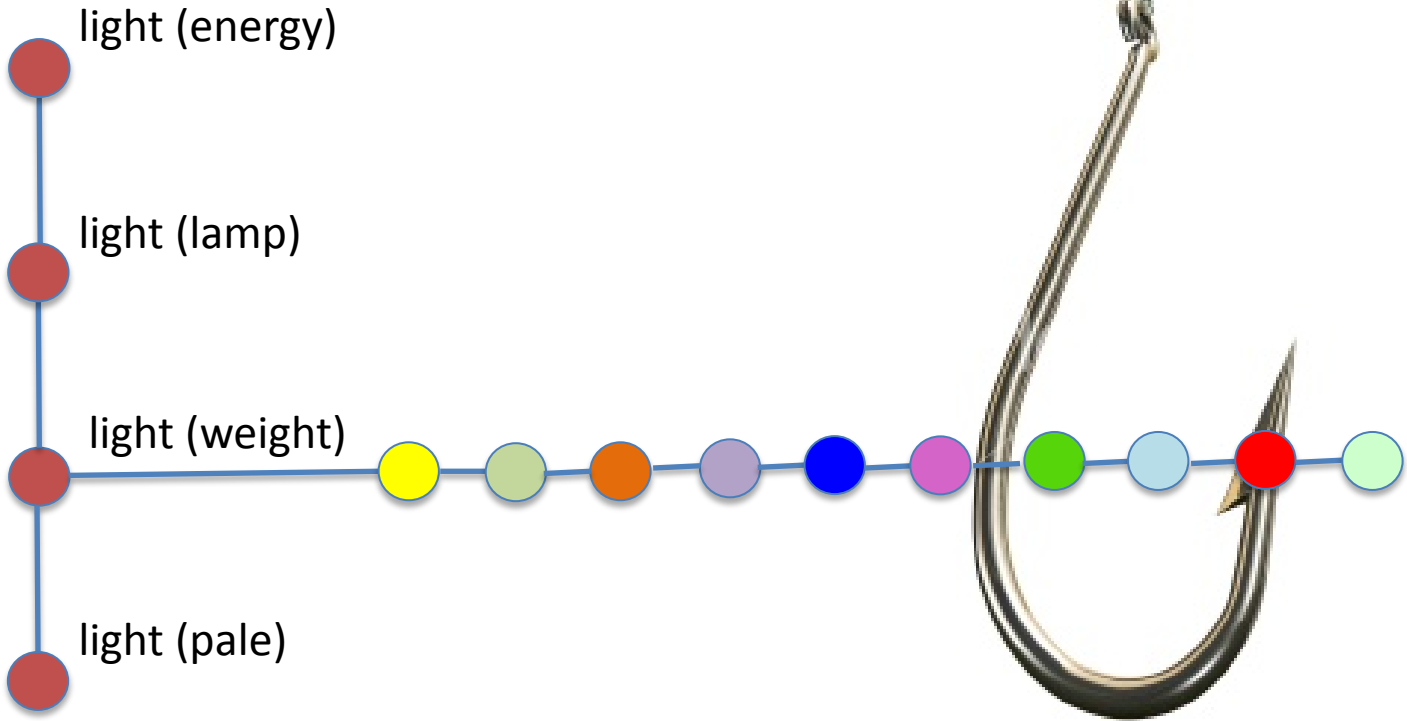
hooks



Resource
Description
Framework

hooks





light (energy)

light (lamp)

light (weight)

light (pale)



1. Project objectives
- 2. Acquiring data**
3. Crowdsourcing innovations



- For people
- For machines

acquiring data



- With people
- With machines

acquiring data



acquiring data

- From existing data
- From experts
- From the crowd

Frairy-n Kapisauan n̄g m̄ga Fraile.
Friction n Pagkuskos.
Friday-n Viernes.
Fried v. imp. & p. p.-Nakaluto; naka-
prito.
Friend-n-Kaulayaw kalaguyo; baibigan;
magkaibigan.

Challenges of existing data
for LRLs:

- Paltry
 - Few data items
 - Few data elements
- Poorly structured
- Inconsistently structured
- No OCR training
- IP restrictions

And then:

- Non-disambiguated

Solution – Merging Engine:

- Get data into consistent fields
- Human review for sense alignment

acquiring data • existing data

Edit Dictionary Term drone

View current Administer Relations Edit current Revisions Devel

Key Information *

Word Forms and Origins

Translations and Concepts

Related terms within the same language

Examples *

Concept

Term (lemma) *	Term Language *	Headword *
<input type="text" value="drone"/>	<input type="text" value="English"/>	<input type="text" value="(?) drone"/>
Part of Speech *	Part Attribute	
<input type="text" value="noun"/>	<input type="text" value="- None -"/>	

Morpheme Information

Collected Morphemes

Plural
<input type="text" value="drones [nid:262211]"/>

Definition

Definition of the term in the language to which it belongs

(?)
An unmanned aerial vehicle (UAV) that is flown for surveillance or military purposes without a human pilot on board.

Definition Source URL

Challenges of language specialists for LRLs:

- Limited # of languages
- Experts have limited time
- Expertise costs money
 - 10,000 entries = 1 year of labor

Expert Tool – Edit Engine:

- Easy
- Consistent
- Comprehensive

acquiring data • experts

Challenges of crowdsourcing for LRLs:

- Eliciting specific data
- Keeping things simple
- Validating data without “ground truth”
 - Minor errors
 - Wrong data
- Malicious users
 - Lots of ways to make mischief
- Sparking participation
- Crowd size
- Maintaining participation

Solutions:

- Fidget Widget
- Facebook Game

acquiring data • the crowd





1. Project objectives
2. Acquiring data
- 3. Crowdsourcing innovations**



Purposes:

- Gather original data
- Manipulate existing data
- Validate existing data
- Validate new data

crowdsourcing



Methods:

- Targeted questions
- Short questions
- Engaging games
- Motivational incentives

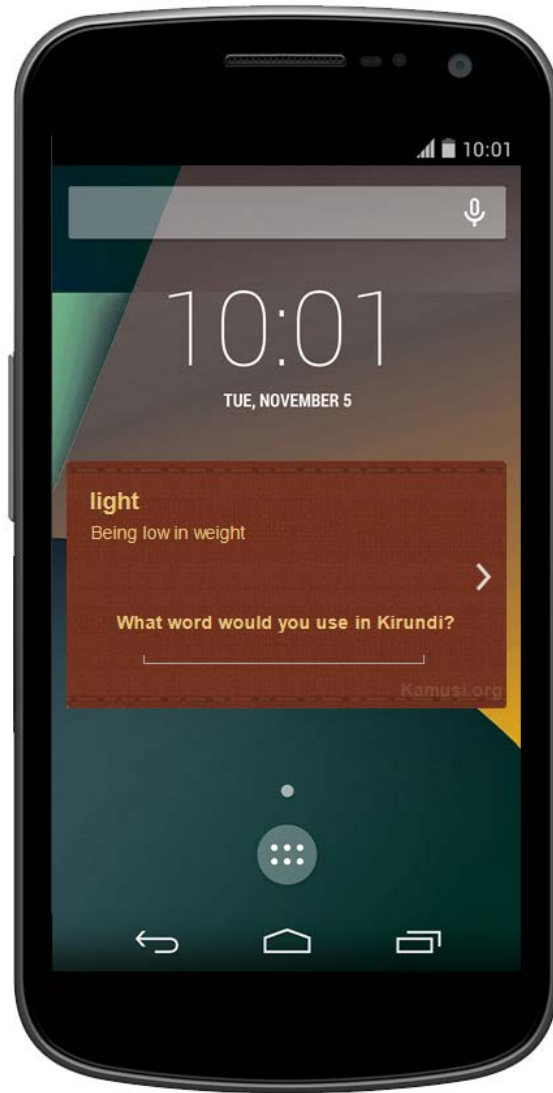
crowdsourcing



Tools:

- Play to Pay
- Fidget Widget (mobile app)
- Facebook Game

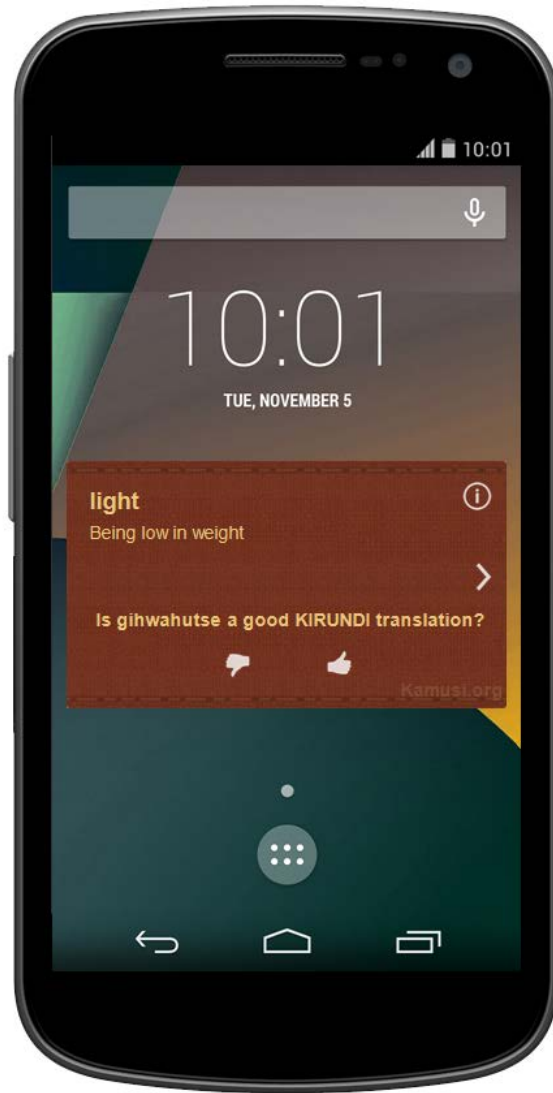
crowdsourcing



Fidget Widget - **Logic**:

- **Single word answers**
- Thumbs up/ down
- Longer answers for advanced users

crowdsourcing



Fidget Widget - Logic:

- Single word answers
- **Thumbs up/ down**
- Longer answers for advanced users

crowdsourcing



Fidget Widget – Logic:

- Single word answers
- Thumbs up/ down
- **Longer answers for advanced users**

crowdsourcing



Fidget Widget - Issues:

- Synchronization
 - Offline in LRL settings
- Security
 - User accounts
- Competing answers
 - Noise vs. signal
- Publication
 - Deciding a quality threshold

crowdsourcing



Fidget Widget – v 2.0:

Talking dictionary

- Voice recording
- Spoken vignettes as definitions
- Spoken translations of definitions

crowdsourcing

Gamification (Fall 2014)

Facebook

- Where the people are
- Where their friends are

The Games

- Pyramid of points for winning answers
- Team points for fastest languages
- Difficulty points for stumpers



crowdsourcing



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE



partners





Contribution analysis

- Good answers – build trust
- Bad answers – change the game
- Malicious answers – remove the user's history

crowdsourcing

languages

5 year goal:
100 languages worldwide

words

5 year goal:
10,000 parallel terms
in each language



data

5 year goal:
100,000,000 searches/ month

people

5 year goal:
1,000,000 registered users



Martin Benjamin, EPFL

Network: <http://www.linkedin.com/in/martinbenjamin>

Write: martin.benjamin@epfl.ch

Meet: <http://meetme.so/kamusi>