

Azaleko sintaxiaren tratamendua  
ikasketa automatikoko tekniken bidez:  
euskarako kateen eta perpausen identifikazioa  
eta bere erabilera koma-zuzentzaile batean

**Doktoregaia:** Bertol Arrieta Kortajarena

**Zuzendariak:** Iñaki Alegria Loinaz

Arantza Diaz de Ilarraza Sanchez

Lengoaia eta Sistema Informatikoak/Lenguajes y Sistemas Informáticos  
Euskal Herriko Unibertsitatea/Universidad del País Vasco

2010eko uztailaren 27a

## Komak eta kateak

*Nork jan du aita?*

## Komak eta kateak

*Nork jan du, aita?*

## Komak eta kateak

*(Nork) (jan du), (aita)?*

## Komak eta perpausak

*Haserretu egin zen emaztea beste batekin ikusi zuenean.*

## Komak eta perpausak

*Haserretu egin zen, emaztea beste batekin ikusi zuenean.*

*Haserretu egin zen emaztea, beste batekin ikusi zuenean.*

## Komak eta perpausak

*(Haserretu egin zen, (emaztea beste batekin ikusi zuenean).)*

*(Haserretu egin zen emaztea, (beste batekin ikusi zuenean).)*

## Tesi honen helburu nagusiak

- 1 Euskarako kate- eta perpaus-identifikatzaile automatikoak sortzea



# Tesi honen helburu nagusiak

- 1 Euskarako kate- eta perpaus-identifikatzaile automatikoak sortzea
- 2 Euskarako koma-zuzentzaile automatikoa garatzea
  - Euskarako komaren erabilera formalizatzea
  - Komaren erabilerak, euskarakoa eta ingelesekoa, konparatzea
  - Kate- eta perpaus-identifikatzaile automatikoen informazioa baliatzea
  - Informazio linguistikoa lortzeko dauzkagun tresnekiko mendekotasuna aztertzea
  - Ebaluazio kualitatiboa egitea

## Bestelako helburuak

- 1 Ikasketa automatikoko teknikak aztertzea
  - Azaleko sintaxian
  - Erroreen detekzioan

## Bestelako helburuak

- 1 Ikasketa automatikoko teknikak aztertzea
  - Azaleko sintaxian
  - Erroreen detekzioan
- 2 Erroreen detekziorako oinarrizko baliabideak garatzea
  - Erroreak etiketatuta dituzten corpusak
  - Euskarako erroreen sailkapena

# Aurkezpenaren eskema

## 1 Testuingurua

# Aurkezpenaren eskema

- 1 Testuingurua
- 2 Kateen eta perpausen identifikazioa

# Aurkezpenaren eskema

- 1 Testuingurua
- 2 Kateen eta perpausen identifikazioa
- 3 Komaren zuzenketa automatikoa

# Aurkezpenaren eskema

- 1 Testuingurua
- 2 Kateen eta perpausen identifikazioa
- 3 Komaren zuzenketa automatikoa
- 4 Ondorioak eta etorkizuneko lanak

# Aurkezpenaren eskema

- 1 Testuingurua
  - Ikasketa automatikoa Hizkuntzaren Prozesamenduan
  - Sintaxiaren tratamendu automatikoa IXA taldean
  - Erroreen detekzio automatikoa IXA taldean
- 2 Kateen eta perpausen identifikazioa
- 3 Komaren zuzenketa automatikoa
- 4 Ondorioak eta etorkizuneko lanak



## Ikasketa automatikoaren funtsa

Ikasketa automatikoa:

Ebatzi beharreko problema sailkapen-ataza batean bilakatzean datza.

# Ikasketa automatikoaren funtsa

## Ikasketa automatikoa:

Ebatzi beharreko problema sailkapen-ataza batean bilakatzean datza.

## Ikasketa gainbegiraturua:

- Ikasi nahi den atazarako corpus etiketatua behar da
- Ikasi nahi den kontzeptuari buruzko adibideekin ikasten du makinak  $\Rightarrow$  adibide berrien portaera iragartzen du
- Bi dimentsioko matrizea: lerroetan, adibideak; zutabeetan, ezaugarriak
- Azken ezaugarria  $\Rightarrow$  ikasi beharreko kontzeptua
- HPan, bereziki sintaxian, hitz bakoitza adibide bat da, eta inguruko hitzen informazioa gehitzeko  $\Rightarrow$  leihoa

# Ikasketa automatikoaren funtsa

Arropa	arropa	IZE	-	B-NP	(S(S*
mota	mota	IZE	-	I-NP	*
guztiak	guzti	DET	ABS	I-NP	*
gustatzen	gustatu	ADI	-	B-VP	*
zaizkidan	izan	ADL	-	I-VP	*
arren	arren	LOT	-	O	*
,	,	PUNT_KOMA	-	O	*S)
nahiago	nahi	IZE	ABS	B-VP	*
ditut	ukan	ADT	-	I-VP	*
soinekoak	soineko	IZE	ABS	B-NP	*
.	.	PUNT_PUNT	-	O	*S)

**Taula:** Perpausen ikasketarako matrize baten adibidea

# Ikasketa automatikoaren funtsa

Arropa	arropa	IZE	-	B-NP	(S(S*
mota	mota	IZE	-	I-NP	*
guztiak	guzti	DET	ABS	I-NP	*
gustatzen	gustatu	ADI	-	B-VP	*
zaizkidan	izan	ADL	-	I-VP	*
arren	arren	LOT	-	O	*
,	,	PUNT_KOMA	-	O	*S)
nahiago	nahi	IZE	ABS	B-VP	*
ditut	ukan	ADT	-	I-VP	*
soinekoak	soineko	IZE	ABS	B-NP	*
.	.	PUNT_PUNT	-	O	*S)

**Taula:** Perpausen ikasketarako matrize baten adibidea

# Ikasketa automatikoaren funtsa

Arropa	arropa	IZE	-	B-NP	(S(S*
mota	mota	IZE	-	I-NP	*
guztiak	guzti	DET	ABS	I-NP	*
gustatzen	gustatu	ADI	-	B-VP	*
zaizkidan	izan	ADL	-	I-VP	*
arren	arren	LOT	-	O	*
,	,	PUNT_KOMA	-	O	*S)
nahiago	nahi	IZE	ABS	B-VP	*
ditut	ukan	ADT	-	I-VP	*
soinekoak	soineko	IZE	ABS	B-NP	*
.	.	PUNT_PUNT	-	O	*S)

Taula: Perpausen ikasketarako matrize baten adibidea

# Ikasketa automatikoaren funtsa

Arropa	arropa	IZE	-	B-NP	(S(S*
mota	mota	IZE	-	I-NP	*
guztiak	guzti	DET	ABS	I-NP	*
gustatzen	gustatu	ADI	-	B-VP	*
zaizkidan	izan	ADL	-	I-VP	*
arren	arren	LOT	-	O	*
,	,	PUNT_KOMA	-	O	*S)
nahiago	nahi	IZE	ABS	B-VP	*
ditut	ukan	ADT	-	I-VP	*
soinekoak	soineko	IZE	ABS	B-NP	*
.	.	PUNT_PUNT	-	O	*S)

Taula: Perpausen ikasketarako matrize baten adibidea

# Ikasketa automatikoaren funtsa

Arropa	arropa	IZE	-	B-NP	(S(S*
mota	mota	IZE	-	I-NP	*
guztiak	guzti	DET	ABS	I-NP	*
gustatzen	gustatu	ADI	-	B-VP	*
<b>zaizkidan</b>	<b>izan</b>	<b>ADL</b>	-	<b>I-VP</b>	<b>*</b>
<b>arren</b>	<b>arren</b>	<b>LOT</b>	-	<b>O</b>	<b>*</b>
,	,	<b>PUNT_KOMA</b>	-	<b>O</b>	<b>*S)</b>
<b>nahiago</b>	<b>nahi</b>	<b>IZE</b>	<b>ABS</b>	<b>B-VP</b>	<b>*</b>
<b>ditut</b>	<b>ukan</b>	<b>ADT</b>	-	<b>I-VP</b>	<b>*</b>
soinekoak	soineko	IZE	ABS	B-NP	*
.	.	PUNT_PUNT	-	O	<b>*S)</b>

Taula: Perpausen ikasketarako matrize baten adibidea

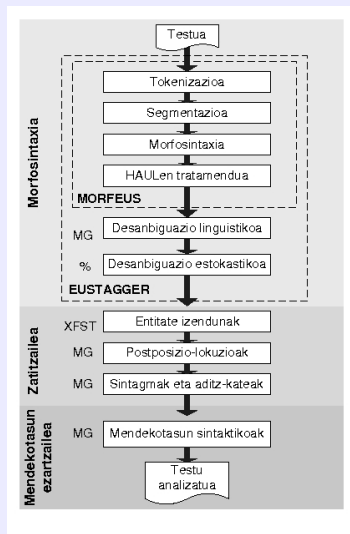
# Oinarrizko ikasketa-algoritmoak

	<b>Estatistikoa</b>	<b>IA</b>	<b>Sinbolikoa</b>	<b>Azpisinbolikoa</b>
<b>Naive Bayes</b>	✓			✓
<b>Erabaki-zuhaitzak</b>		✓	✓	
<b>Pertzeptroiak</b>		✓		✓
<b>SVM</b>		✓		✓

**Taula:** Tesi-lan honetan erabilitako oinarrizko ikasketa-eskemen sailkapena



# Sintaxiaren tratamendu automatikoa IXA taldean



## Euskarako erroren detekzioa IXA taldean

Datak, postposizio-lokuzioak eta komunztadura (Ornoz, 2009)

**\*Zentral nuklearrak *zakar erradiaktiboa* eratzen dute.**

## Euskarako erroreen detekzioa IXA taldean

Datak, postposizio-lokuzioak eta komunztadura (Ornoz, 2009)

**\*Zentral nuklearrak** *zakar erradiaktiboa* **eratzten dute.**

Determinatzaileak (Uria, 2009)

**\*Amak esan dit** **pinguino bezala** *nabilela.*

## Euskarako erroren detekzioa IXA taldean

Datak, postposizio-lokuzioak eta komunztadura (Ornoz, 2009)

**\*Zentral nuklearrak zakar erradiaktiboa eratzen dute.**

Determinatzaileak (Uria, 2009)

**\*Amak esan dit pinguino bezala nabilela.**

Ikasketa automatiko bidez:

- Proba txiki bat determinatzaile- eta komunztadura-akatsak detektatzeko: corpus txikia, errore gutxi
- Corpus egokia behar
- Komen ikasketa: corpora etiketatzeko beharrik ez

# Aurkezpenaren eskema

- 1 Testuingurua
- 2 Kateen eta perpausen identifikazioa
  - Testuingurua
  - FR-Perceptron
  - Esperimentuen prestaketa
  - Kateen identifikazio automatikoa
  - Perpausen identifikazio automatikoa
- 3 Komaren zuzenketa automatikoa
- 4 Ondorioak eta etorkizuneko lanak

## Hitz multzoak: definizio formal bat (Carreras, 2005)

$$\mathcal{P} = \{(s, e)_k \mid 1 \leq s \leq e, k \in \mathcal{K}\}$$

# Hitz multzoak: definizio formal bat (Carreras, 2005)

$$\mathcal{P} = \{(s, e)_k \mid 1 \leq s \leq e, k \in \mathcal{K}\}$$

$p_1 = (s_1, e_1)$  eta  $p_2 = (s_2, e_2)$  hitz multzoak **gainjartzen** dira ( $p_1 \sim p_2$ ), baldin:

$$s_1 < s_2 \leq e_1 < e_2 \text{ edo } s_2 < s_1 \leq e_2 < e_1$$

Grafikoki:  $(s_1 \quad (s_2 \quad )e_1 \quad )e_2$

# Hitz multzoak: definizio formal bat (Carreras, 2005)

$$\mathcal{P} = \{(s, e)_k \mid 1 \leq s \leq e, k \in \mathcal{K}\}$$

$p_1 = (s_1, e_1)$  eta  $p_2 = (s_2, e_2)$  hitz multzoak **gainjartzen** dira ( $p_1 \sim p_2$ ), baldin:

$$s_1 < s_2 \leq e_1 < e_2 \text{ edo } s_2 < s_1 \leq e_2 < e_1$$

Grafikoki:  $(s_1 \quad (s_2 \quad )e_1 \quad )e_2$

$p_2$  hitz multzoak  $p_1$  hitz multzoa **bere barnean hartzen** du ( $p_1 \preceq p_2$ ), baldin:

$$s_2 \leq s_1 \leq e_1 \leq e_2$$

Grafikoki:  $(s_2 \quad (s_1 \quad )e_1 \quad )e_2$



# Kateak eta perpausak hitz multzo gisa

Kateen soluzio-espazioa (Abney, 1991)

$$\mathcal{Y} = \{y \subseteq \mathcal{P} \mid \forall p_1, p_2 \in y (p_1 \not\sim p_2 \wedge p_1 \not\neq p_2)\}$$

*(Lehendakariak) (bere erabakia) (berretsi zuen).*

# Kateak eta perpausak hitz multzo gisa

## Kateen soluzio-espazioa (Abney, 1991)

$$\mathcal{Y} = \{y \subseteq \mathcal{P} \mid \forall p_1, p_2 \in y (p_1 \not\sim p_2 \wedge p_1 \not\neq p_2)\}$$

*(Lehendakariak) (bere erabakia) (berretsi zuen).*

## Perpausen soluzio-espazioa:

$$\mathcal{Y} = \{y \subseteq \mathcal{P} \mid \forall p_1, p_2 \in y (p_1 \not\sim p_2)\}$$

*(Zubietarrak kezkatuta daude, (planak euren herria desitxuratuko duelako).)*

## Kateak: zenbait adibide

Ingelesezt: lana banatzearen onurak

*(The deficit) (will narrow) (to) (only 1.8 billion) (in) (September).*

*(The deficit) (will narrow) (to only 1.8 billion) (in September).*

## Kateak: zenbait adibide

Ingelesez: lana banatzearen onurak

*(The deficit) (will narrow) (to) (only 1.8 billion) (in) (September).*  
*(The deficit) (will narrow) (to only 1.8 billion) (in September).*

Euskarako kateak: barne-hartzea?

*(Ez dut (horretan) sakondu nahi).*  
*(Ez dut) (horretan) (sakondu nahi).*

## Perpausak: zenbait adibide

### Perpausen barnean perpausak:

*((Hiriko agintari militarrek aditzera eman dutenez,) suziri batek eragin zuen leherketa.)*

*((Urteko helburu nagusia (bakea lortzea) izango dela) adierazi du.)*

## Perpausak: zenbait adibide

### Perpausen barnean perpausak:

*((Hiriko agintari militarrek aditzera eman dutenez,) suziri batek eragin zuen leherketa.)*

*((Urteko helburu nagusia (bakea lortzea) izango dela) adierazi du.)*

### Anbiguotasuna

*(Haserretu egin zen (emaztea beste batekin ikusi zuenean).)*

*(Haserretu egin zen emaztea (beste batekin ikusi zuenean).)*

## Ingeleseko kateen identifikazioa

	<b>Teknika</b>	$F_1$
<b>(Lee eta Wu, 2007)</b>	<i>SVM</i>	94,22
<b>(Zhang et al., 2002)</b>	<i>Winnow</i> orokortua	94,13
<b>(Shen eta Sarkar, 2005)</b>	<i>Voted HMM</i>	94,01
<b>(Kudo eta Matsumoto, 2001)</b>	<i>SVM</i>	93,91
<b>(Carreras et al., 2005)</b>	<i>FR-perceptron</i>	93,74
<b>(Molina eta Pla, 2002)</b>	<i>Spec. HMM</i>	93,25
<b>Oinarrizko neurria (baseline-a)</b>	kategoriaren usuena	77,07

**Taula:** Ingeleseko kate-identifikatzaile onenak, *CoNLL 2000*ko baldintzetan

## Ingeleseko kateen identifikazioa

	<b>Teknika</b>	$F_1$
(Lee eta Wu, 2007)	<i>SVM</i>	94,22
(Zhang et al., 2002)	<i>Winnow</i> orokortua	94,13
(Shen eta Sarkar, 2005)	<i>Voted HMM</i>	94,01
(Kudo eta Matsumoto, 2001)	<i>SVM</i>	93,91
(Carreras et al., 2005)	<i>FR-perceptron</i>	93,74
(Molina eta Pla, 2002)	<i>Spec. HMM</i>	93,25
<b>Oinarrizko neurria (baseline-a)</b>	<b>kategoriaren usuena</b>	<b>77,07</b>

**Taula:** Ingeleseko kate-identifikatzaile onenak, *CoNLL 2000*ko baldintzetan



# Ingeleseko kateen identifikazioa

	<b>Teknika</b>	$F_1$
(Lee eta Wu, 2007)	<i>SVM</i>	94,22
(Zhang et al., 2002)	<i>Winnow orokortua</i>	94,13
(Shen eta Sarkar, 2005)	<i>Voted HMM</i>	94,01
(Kudo eta Matsumoto, 2001)	<i>SVM</i>	93,91
(Carreras et al., 2005)	<b>FR-perceptron</b>	<b>93,74</b>
(Molina eta Pla, 2002)	<i>Spec. HMM</i>	93,25
<b>Oinarrizko neurria (baseline-a)</b>	kategoriaren usuena	77,07

**Taula:** Ingeleseko kate-identifikatzaile onenak, *CoNLL 2000*ko baldintzetan

## Ingeleseko perpausen identifikazioa

	<b>Teknika</b>	$F_1$
<b>(Ram eta Devi, 2008)</b>	<i>CRF + erregela ling.</i>	89,04
<b>(Carreras et al., 2005)</b>	<i>FR-Perceptron</i>	85,03
<b>(Nguyen et al., 2009)</b>	<i>joint-CRF</i>	84,66
<b>(Carreras eta Màrquez, 2001)</b>	<i>Boosting</i>	81,73
<b>(Molina eta Pla, 2001)</b>	<i>HMM</i>	70,68
<b>(Tjong Kim Sang, 2001)</b>	<i>MBL</i>	70,58
<b>(Patrick eta Goyal, 2001)</b>	<i>Erabaki-grafoak</i>	68,85
<b>Oinarrizko neurria (baseline-a)</b>	Esaldiak mugatu	47,71

**Taula:** Ingeleseko perpaus-identifikatzaile onenak, *CoNLL 2001*eko baldintzetan

## Ingeleseko perpausen identifikazioa

	<b>Teknika</b>	$F_1$
<b>(Ram eta Devi, 2008)</b>	<i>CRF + erregela ling.</i>	89,04
<b>(Carreras et al., 2005)</b>	<i>FR-Perceptron</i>	85,03
<b>(Nguyen et al., 2009)</b>	<i>joint-CRF</i>	84,66
<b>(Carreras eta Màrquez, 2001)</b>	<i>Boosting</i>	81,73
<b>(Molina eta Pla, 2001)</b>	<i>HMM</i>	70,68
<b>(Tjong Kim Sang, 2001)</b>	<i>MBL</i>	70,58
<b>(Patrick eta Goyal, 2001)</b>	<i>Erabaki-grafoak</i>	68,85
<b>Oinarrizko neurria (baseline-a)</b>	<b>Esaldiak mugatu</b>	<b>47,71</b>

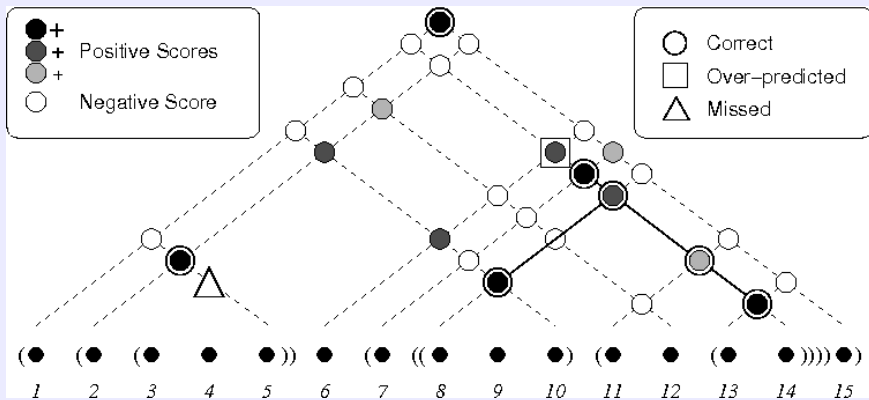
Taula: Ingeleseko perpaus-identifikatzaile onenak, CoNLL 2001eko baldintzetan

## Ingeleseko perpausen identifikazioa

	<b>Teknika</b>	$F_1$
(Ram eta Devi, 2008)	<i>CRF + erregela ling.</i>	89,04
(Carreras et al., 2005)	<b>FR-Perceptron</b>	<b>85,03</b>
(Nguyen et al., 2009)	<i>joint-CRF</i>	84,66
(Carreras eta Màrquez, 2001)	<i>Boosting</i>	81,73
(Molina eta Pla, 2001)	<i>HMM</i>	70,68
(Tjong Kim Sang, 2001)	<i>MBL</i>	70,58
(Patrick eta Goyal, 2001)	<i>Erabaki-grafoak</i>	68,85
Oinarrizko neurria (baseline-a)	Esaldiak mugatu	47,71

Taula: Ingeleseko perpaus-identifikatzaile onenak, CoNLL 2001eko baldintzetan

# FR-Perceptron algoritmoa: adibide bat (Carreras, 2005)



## Esperimentuen prestaketa (I)

- Corputa:

	<b>Guztira</b>	<b>Ikasketa-corputa</b>	<b>Garapen-corputa</b>	<b>Test-corputa</b>
<i>EPEC</i>	150.128	104.956	22.548	22.624

**Taula:** Kateen eta perpausen identifikaziorako corputa

- Ebaluaziorako neurriak:

Doitasuna, estaldura eta  $F_1$  neurria.

$$F_1 = \frac{2 * Doitasuna * Estaldura}{(Doitasuna + Estaldura)}$$

## Esperimentuen prestaketa (II): oinarrizko neurriak

		<b>Doitasuna</b>	<b>Estaldura</b>	$F_1$ <b>neurria</b>
<b>Euskarako oinarrizko neurriak</b>	<b>Kateak</b>	45,76	60,21	52,00
<b>Ingeleseko oinarrizko neurriak</b>	<b>Kateak</b>	72,58	82,14	77,07

Taula: **Kateen** identifikazioko oinarrizko neurriak: euskara vs ingelesa

## Esperimentuen prestaketa (II): oinarrizko neurriak

		Doitasuna	Estaldura	$F_1$ neurria
<b>Euskarako oinarrizko neurriak</b>	<b>Kateak</b>	45,76	60,21	<b>52,00</b>
<b>Ingeleseko oinarrizko neurriak</b>	<b>Kateak</b>	72,58	82,14	<b>77,07</b>

Taula: **Kateen** identifikazioko oinarrizko neurriak: euskara vs ingelesa



## Esperimentuen prestaketa (II): oinarrizko neurriak

		Doitasuna	Estaldura	$F_1$ neurria
Euskarako oinarrizko neurriak	Sint.	31,50	47,08	37,74
	ADK	77,65	80,63	79,11
	Kateak	45,76	60,21	52,00
Ingeleseko oinarrizko neurriak	Kateak	72,58	82,14	77,07

Taula: Kateen identifikazioko oinarrizko neurriak: euskara vs ingelesa

## Esperimentuen prestaketa (II): oinarrizko neurriak

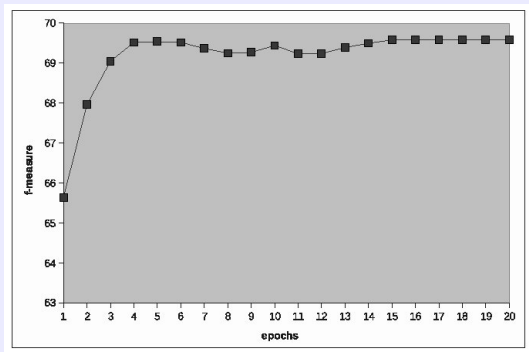
	<b>Doitasuna</b>	<b>Estaldura</b>	<b><math>F_1</math> neurria</b>
<b>Euskarako oinarrizko neurriak</b>	91,41	33,27	<b>48,79</b>
<b>Ingeleseko oinarrizko neurriak</b>	98,44	31,48	<b>47,71</b>

Taula: **Perpausen** identifikazioko oinarrizko neurriak: euskara vs ingelesa

# Kateen identifikazio automatikoa

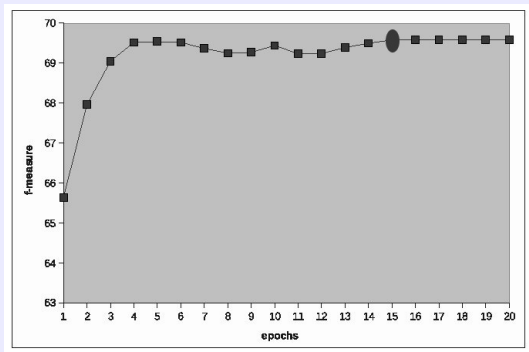
- Epoch zenbakiaren eragina
- Lehen probak oinarrizko ezaugarriekin
- Ezaugarri linguistikoak gehituz
- Kateen identifikazioa IA eta erregelak konbinatuz
- Ikasketa-corpora handituz
- Emaizta idealak eskuz desanbiguatutako informazioarekin
- Azken emaitzak test-corpusean

## Epoch zenbakiaren eragina



**Irudia:** *Epoch-zenbakiaren eragina FR-Perceptron bidezko ikasketa automatikoan*

## Epoch zenbakiaren eragina



**Irudia:** *Epoch-zenbakiaren eragina FR-Perceptron bidezko ikasketa automatikoan*

## Lehen probak oinarritzko ezaugarriekin

		Doitasuna	Estaldura	$F_1$ neurria
Oinarritzko neurriak (baseline-a)	Kateak	45,76	60,21	52,00
FR-Perceptron oinarritzko ezaugarriekin	Kateak	70,85	68,36	69,58

Taula: *FR-Perceptron* algoritmoan oinarritutako kate-identifikatzailearen emaitzak, oinarritzko ezaugarriak erabilita: **hitza eta kategoria**

## Lehen probak oinarrizko ezaugarriekin

		Doitasuna	Estaldura	$F_1$ neurria
Oinarrizko neurriak (baseline-a)	Kateak	45,76	60,21	<b>52,00</b>
FR-Perceptron oinarrizko ezaugarriekin	Kateak	70,85	68,36	<b>69,58</b>

*Taula:* FR-Perceptron algoritmoan oinarritutako kate-identifikatzailearen emaitzak, oinarrizko ezaugarriak erabilita: **hitza eta kategoria**

## Lehen probak oinarrizko ezaugarriekin

		Doitasuna	Estaldura	$F_1$ neurria
Oinarrizko neurriak (baseline-a)	Sint.	31,50	47,08	37,74
	ADK	77,65	80,63	79,11
	Kateak	45,76	60,21	52,00
FR-Perceptron oinarrizko ezaugarriekin	Sint.	60,54	55,49	57,90
	ADK	84,98	88,40	86,66
	Kateak	70,85	68,36	69,58

*Taula:* FR-Perceptron algoritmoan oinarritutako kate-identifikatzailearen emaitzak, oinarrizko ezaugarriak erabilita: **hitza eta kategoria**



## Lehen probak oinarritzko ezaugarriekin

		Doitasuna	Estaldura	$F_1$ neurria
Oinarritzko neurriak (baseline-a)	Sint.	31,50	47,08	<b>37,74</b>
	ADK	77,65	80,63	<b>79,11</b>
	Kateak	45,76	60,21	52,00
FR-Perceptron oinarritzko ezaugarriekin	Sint.	60,54	55,49	<b>57,90</b>
	ADK	84,98	88,40	<b>86,66</b>
	Kateak	70,85	68,36	69,58

*Taula:* *FR-Perceptron* algoritmoan oinarritutako kate-identifikatzailearen emaitzak, *oinarritzko ezaugarriak* erabilita: **hitza eta kategoria**

## Ezaugarri linguistikoak gehituz

	<b>Doitasuna</b>	<b>Estaldura</b>	<b><math>F_1</math> neurria</b>
<b>oin</b>	70,85	68,36	69,58
<b>oin + ak</b>	70,78	69,15	69,96
<b>oin + dek</b>	77,42	80,17	78,77
<b>oin + l</b>	71,53	68,91	70,20
<b>oin + men</b>	70,77	68,48	69,61
<b>oin + ak + dek + l + men</b>	77,67	79,70	78,67

**Taula:** Ezaugarri linguistikoak gehituz (ak: azpikategoria; dek: deklinabidea; l: lema; men: mendeko perpaus mota)

## Ezaugarri linguistikoak gehituz

	Doitasuna	Estaldura	$F_1$ neurria
oin	70,85	68,36	69,58
oin + ak	70,78	69,15	69,96
oin + dek	77,42	80,17	<b>78,77</b>
oin + l	71,53	68,91	70,20
oin + men	70,77	68,48	69,61
oin + ak + dek + l + men	77,67	79,70	78,67

**Taula:** Ezaugarri linguistikoak gehituz (ak: azpikategoria; dek: deklinabidea; l: lema; men: mendeko perpaus mota)

## Deklinabidearen garrantzia sintagmetan

	<b>Sint. <math>F_1</math> neurria</b>	<b>ADK <math>F_1</math> neurria</b>	<b><math>F_1</math> neurria</b>
<b>Oin</b>	57,90	86,66	69,58
<b>Oin + dek</b>	73,28	87,33	78,77

**Taula:** Deklinabidearen eraginaren konparazioa sintagmetan eta aditz-kateetan

## Deklinabidearen garrantzia sintagmetan

	Sint. $F_1$ neurria	ADK $F_1$ neurria	$F_1$ neurria
<b>Oin</b>	57,90	<b>86,66</b>	69,58
<b>Oin + dek</b>	73,28	<b>87,33</b>	78,77

**Taula:** Deklinabidearen eraginaren konparazioa sintagmetan eta aditz-kateetan

## Deklinabidearen garrantzia sintagmetan

	Sint. $F_1$ neurria	ADK $F_1$ neurria	$F_1$ neurria
Oin	<b>57,90</b>	86,66	69,58
Oin + dek	<b>73,28</b>	87,33	78,77

**Taula:** Deklinabidearen eraginaren konparazioa sintagmetan eta aditz-kateetan

## Kateen identifikazioa IA eta erregelak konbinatuz

	<b>Doitasuna</b>	<b>Estaldura</b>	<b><math>F_1</math> neurria</b>
<b>Erreg</b>	50,06	52,98	51,48
<b>FR-Perceptron oin</b>	70,85	68,36	69,58
<b>FR-Perceptron oin + erreg</b>	75,51	77,61	76,54
<b>FR-Perceptron oin + dek</b>	77,42	80,17	78,77
<b>FR-Perceptron oin + dek + erreg</b>	77,68	80,80	79,21

**Taula:** Erregelaren informazioa gehitzea, pilaratzearen bidez

## Kateen identifikazioa IA eta erregelak konbinatuz

	Doitasuna	Estaldura	$F_1$ neurria
Erreg	50,06	52,98	51,48
FR-Perceptron oin	70,85	68,36	<b>69,58</b>
FR-Perceptron oin + erreg	75,51	77,61	<b>76,54</b>
FR-Perceptron oin + dek	77,42	80,17	78,77
FR-Perceptron oin + dek + erreg	77,68	80,80	79,21

Taula: Erregelaren informazioa gehitzea, pilaratzearen bidez



## Kateen identifikazioa IA eta erregelak konbinatuz

	Doitasuna	Estaldura	$F_1$ neurria
Erreg	50,06	52,98	51,48
FR-Perceptron oin	70,85	68,36	69,58
FR-Perceptron oin + erreg	75,51	77,61	76,54
FR-Perceptron oin + dek	77,42	80,17	<b>78,77</b>
FR-Perceptron oin + dek + erreg	77,68	80,80	<b>79,21</b>

Taula: Erregelaren informazioa gehitzea, pilaratzearen bidez

## Ikasketa-corpusa handituz

	<b>Doitasuna</b>	<b>Estaldura</b>	<b><math>F_1</math> neurria</b>
<b>Ikasketa-corpusaren % 25arekin</b>	77,68	80,80	79,21
<b>Ikasketa-corpusaren % 50arekin</b>	79,46	83,06	81,22
<b>Ikasketa-corpusaren % 100arekin</b>	81,09	84,24	82,64

Taula: Ikasketa-corpusaren tamainaren arabera (% 100 = 104.956 token)

# Ikasketa-corpusa handituz

	Doitasuna	Estaldura	$F_1$ neurria
<b>Ikasketa-corpusaren % 25arekin</b>	77,68	80,80	<b>79,21</b>
<b>Ikasketa-corpusaren % 50arekin</b>	79,46	83,06	<b>81,22</b>
<b>Ikasketa-corpusaren % 100arekin</b>	81,09	84,24	82,64

Taula: Ikasketa-corpusaren tamainaren arabera (% 100 = 104.956 token)

## Ikasketa-corpusa handituz

	Doitasuna	Estaldura	$F_1$ neurria
Ikasketa-corpusaren % 25arekin	77,68	80,80	79,21
Ikasketa-corpusaren % 50arekin	79,46	83,06	<b>81,22</b>
Ikasketa-corpusaren % 100arekin	81,09	84,24	<b>82,64</b>

Taula: Ikasketa-corpusaren tamainaren arabera (% 100 = 104.956 token)

## Emaitza idealak eskuz desanbiguatutako informazioarekin

Corpusa nola desanbiguatua	Doit.	Est.	Sint. $F_1$	ADK $F_1$	$F_1$
<i>Automatikoki</i>	81,09	84,24	78,00	89,84	82,64
<i>Eskuz</i>	89,61	91,46	87,85	94,52	90,52

**Taula:** Erabilitako informazio linguistikoaren kalitatearen arabera

## Emaitza idealak eskuz desanbiguatutako informazioarekin

Corpusa nola desanbiguatua	Doit.	Est.	Sint. $F_1$	ADK $F_1$	$F_1$
<i>Automatikoki</i>	81,09	84,24	78,00	89,84	<b>82,64</b>
<i>Eskuz</i>	89,61	91,46	87,85	94,52	<b>90,52</b>

Taula: Erabilitako informazio linguistikoaren kalitatearen arabera

## Emitza idealak eskuz desanbiguatutako informazioarekin

Corpusa nola desanbiguatua	Doit.	Est.	Sint. $F_1$	ADK $F_1$	$F_1$
<i>Automatikoki</i>	81,09	84,24	<b>78,00</b>	89,84	<b>82,64</b>
<i>Eskuz</i>	89,61	91,46	<b>87,85</b>	94,52	<b>90,52</b>

Taula: Erabilitako informazio linguistikoaren kalitatearen arabera

## Emaitza idealak eskuz desanbiguatutako informazioarekin

Corpusa nola desanbiguatua	Doit.	Est.	Sint. $F_1$	ADK $F_1$	$F_1$
<i>Automatikoki</i>	81,09	84,24	<b>78,00</b>	89,84	<b>82,64</b>
<i>Eskuz</i>	89,61	91,46	<b>87,85</b>	94,52	<b>90,52</b>

Taula: Erabilitako informazio linguistikoaren kalitatearen arabera

Desanbiguzio maila	$F_1$
1. mailan	95,92
2. mailan (ak)	95,42
3. mailan (ak + dek)	<b>90,36</b>

Taula: *Eustagger*-en emaitzak desanbiguzio mailaren arabera (Ezeiza, 2002)



## Azken emaitzak test-corpusean

Proba-corpusa	Doit.	Est.	$F_1$
Garapen-corpusa	81,09	84,24	<b>82,64</b>
Test-corpusa	81,35	85,07	<b>83,17</b>

Taula: Garapen-corpuseko emaitzak vs test-corpuseko emaitzak

# Perpausen identifikazio automatikoa

- Lehen probak oinarrizko ezaugarriekin
- Ezaugarri linguistikoak gehituz
- Kateen identifikazioa IA eta erregelak konbinatuz
- Ikasketa-corpora handituz
- Emaidza idealak eskuz desanbiguatutako informazioarekin
- Azken emaitzak test-corpusean

## Lehen probak oinarrizko ezaugarriekin

	Doitasuna	Estaldura	$F_1$ neurria
<b>Oinarrizko neurriak (baseline-a)</b>	91,41	33,27	48,79
<b>FR-Perc. oinarrizko ezaugarriekin</b>	76,72	64,34	<b>69,99</b>

**Taula:** Oinarrizko ezaugarriekin: **hitz, kategoria eta kateei buruzko informazioa**

## Ezaugarri linguistikoak gehituz

	<b>Doitasuna</b>	<b>Estaldura</b>	<b><math>F_1</math> neurria</b>
<b>oin</b>	76,72	64,34	69,99
<b>oin + ak</b>	76,52	66,49	71,15
<b>oin + dek</b>	76,41	66,10	70,89
<b>oin + l</b>	75,66	67,58	71,39
<b>oin + men</b>	76,39	65,68	70,63
<b>oin + ak + dek + l + men</b>	76,82	69,94	73,22

**Taula:** Ezaugarri linguistikoak gehituz (ak: azpikategoria; dek: deklinabidea; l: lema; men: mendeko perpaus mota)

## Ezaugarri linguistikoak gehituz

	Doitasuna	Estaldura	$F_1$ neurria
oin	76,72	64,34	69,99
oin + ak	76,52	66,49	71,15
oin + dek	76,41	66,10	70,89
oin + l	75,66	67,58	71,39
oin + men	76,39	65,68	70,63
oin + ak + dek + l + men	76,82	69,94	<b>73,22</b>

**Taulara:** Ezaugarri linguistikoak gehituz (ak: azpikategoria; dek: deklinabidea; l: lema; men: mendeko perpaus mota)

## Perpausen identifikazioa IA eta erregelak konbinatuz

	Doitasuna	Estaldura	$F_1$
<b>erreg</b>	50,84	48,63	49,71
<b>oin + ak + dek + l + men</b>	76,82	69,94	73,22
<b>oin + ak + dek + l + men + erreg</b>	78,03	71,35	74,54

**Taula:** Erregelaren informazioa gehitzea, pilaratzearen bidez

## Perpausen identifikazioa IA eta erregelak konbinatuz

	Doitasuna	Estaldura	$F_1$
erreg	50,84	48,63	49,71
oin + ak + dek + l + men	76,82	69,94	73,22
oin + ak + dek + l + men + erreg	78,03	71,35	<b>74,54</b>

Taula: Erregelaren informazioa gehitzea, pilaratzearen bidez

# Ikasketa-corpusa handituz

	<b>Doitasuna</b>	<b>Estaldura</b>	$F_1$
<b>Ikasketa-corpusaren % 25arekin</b>	78,03	71,35	74,54
<b>Ikasketa-corpusaren % 50arekin</b>	78,70	74,21	76,39
<b>Ikasketa-corpusaren % 100arekin</b>	80,13	76,18	78,11

**Taula:** Corpusaren tamainaren arabera (% 100 = 104.956 token)



# Ikasketa-corpora handituz

	Doitasuna	Estaldura	$F_1$
Ikasketa-corporaren % 25arekin	78,03	71,35	<b>74,54</b>
Ikasketa-corporaren % 50arekin	78,70	74,21	<b>76,39</b>
Ikasketa-corporaren % 100arekin	80,13	76,18	78,11

Taula: Corpusaren tamainaren arabera (% 100 = 104.956 token)

# Ikasketa-corpora handituz

	Doitasuna	Estaldura	$F_1$
Ikasketa-corporaren % 25arekin	78,03	71,35	74,54
Ikasketa-corporaren % 50arekin	78,70	74,21	<b>76,39</b>
Ikasketa-corporaren % 100arekin	80,13	76,18	<b>78,11</b>

Taula: Corpusaren tamainaren arabera (% 100 = 104.956 token)

## Emitza idealak eskuz desanbiguatutako informazioarekin

<b>Corpusa nola desanbiguatua</b>	<b>Doit.</b>	<b>Est.</b>	<b><math>F_1</math></b>
<i>Automatikoki</i>	80,13	76,18	<b>78,11</b>
<i>Eskuz</i>	80,81	76,11	<b>78,39</b>

Taula: Informazio linguistikoaren kalitatearen arabera

## Azken emaitzak test-corpusean

Proba-corpusa	Doit.	Est.	$F_1$
Garapen-corpusa	80,13	76,18	<b>78,11</b>
Test-corpusa	79,22	75,36	<b>77,24</b>

Taula: Garapen-corpuseko emaitzak vs test-corpuseko emaitzak

# Aurkezpenaren eskema

- 1 Testuingurua
- 2 Kateen eta perpausen identifikazioa
- 3 **Komaren zuzenketa automatikoa**
  - Puntuazioa HPan
  - Komaren erabilera: azterketa linguistikoa
  - Komen zuzenketa hizkuntzaren ezagutzan oinarrituta
  - Komen zuzenketa ikasketa automatikoan oinarrituta
- 4 Ondorioak eta etorkizuneko lanak

## Puntuazioa HPan

- Puntuazioaren garrantzia HPan
  - Komak esaldiaren sintaxian daukan garrantzia (Nunberg, 1990)
  - Puntuazioaren eragina analizatzaile sintaktiko batean (Briscoe eta Carroll, 1995)
  - *Workshop on punctuation in computational linguistics (1996)*

## Puntuazioa HPan

- Puntuazioaren garrantzia HPan
  - Komak esaldiaren sintaxian daukan garrantzia (Nunberg, 1990)
  - Puntuazioaren eragina analizatzaile sintaktiko batean (Briscoe eta Carroll, 1995)
  - *Workshop on punctuation in computational linguistics (1996)*
- Lan esanguratsuenak
  - Komaren funtzioen sailkapena (Bayraktar et al., 1998)
  - Danierako koma-zuzentzailea (Hardt, 2001)
  - Komen rol sintaktikoen esleipena (Delden eta Gomez, 2002)
  - Komen berreskurapena ahotsaren ezagutzarako (Shieber eta Tao, 2003)
  - Puntuazioaren eta kasuaren berreskurapena (Baldwin eta Joseph, 2009)
  - Txekierako puntuazio-akatsak detektatzeko (Jakubicek eta Horak, 2010)

## Puntuazioa HPan

- Puntuazioaren garrantzia HPan
  - Komak esaldiaren sintaxian daukan garrantzia (Nunberg, 1990)
  - Puntuazioaren eragina analizatzaile sintaktiko batean (Briscoe eta Carroll, 1995)
  - *Workshop on punctuation in computational linguistics (1996)*
- Lan esanguratsuenak
  - Komaren funtzioen sailkapena (Bayraktar et al., 1998)
  - Danierako koma-zuzentzailea (Hardt, 2001)
  - Komen rol sintaktikoen esleipena (Delden eta Gomez, 2002)
  - **Komen berreskurapena ahotsaren ezagutzarako (Shieber eta Tao, 2003)**
  - Puntuazioaren eta kasuaren berreskurapena (Baldwin eta Joseph, 2009)
  - Txekierako puntuazio-akatsak detektatzeko (Jakubicek eta Horak, 2010)



## Puntuazio-zuzentzaile bat garatzeko arazoak

- Puntuazio-arauak hizkuntza askotan ez daude guztiz zehaztuta
- Puntuazio-marka fidagarriak (puntuak, galdera-marka eta harridura-marka) vs ez-fidagarriak
- Komaren kasu berezia

## Komaren erabilera: azterketa linguistikoa

- Aditz nagusiaren aurreko mintzagaiaren ondoren koma jartzea gomendatzen da (1), mintzagaia subjektua ez denean betiere (2). Aitzitik, behar-beharrezkoa da batzuetan, anbiguotasuna ekiditeko (3). Mintzagaia edo galdegaia (bietako bat, gutxienez) mendeko perpausa bada, mintzagaiaren ondoren koma jarri behar da derrigor (4,5) (Garzia, 1997).
  - 1 **Azkenean**, *zurekin joango naiz.*
  - 2 **Aita** *atzo iritsi zen.*
  - 3 **Batzuetan**, *irabazteko gogor jokatzera beharrezkoa da.*
  - 4 **Euria ari zuenez**, *etxean geratzea erabaki genuen.*
  - 5 **Gaur iritsiko zela esan zuen arren**, *atzo iritsi zen azkenean.*

## Komaren erabilera: azterketa linguistikoa

- Aditz nagusiaren aurreko mintzagaiaren ondoren koma jartzea gomendatzen da (1), mintzagaia subjektua ez denean betiere (2). Aitzitik, behar-beharrezkoa da batzuetan, anbiguotasuna ekiditeko (3). **Mintzagaia edo galdegaia (bietako bat, gutxienez) mendeko perpausa bada, mintzagaiaren ondoren koma jarri behar da derrigor (4,5)** (Garzia, 1997).
  - ① **Azkenean**, *zurekin joango naiz.*
  - ② **Aita** *atzo iritsi zen.*
  - ③ **Batzuetan**, *irabazteko gogor jokatzea beharrezkoa da.*
  - ④ **Euria ari zuenez**, *etxean geratzea erabaki genuen.*
  - ⑤ **Gaur iritsiko zela esan zuen arren**, *atzo iritsi zen azkenean.*

## Komaren erabileraren konparaketa: euskara eta ingelesa

Euskarakoak	Ingelesekoak
Esaldi koordinatuak lotzean, juntagailuaren aurretik koma	Enumerazio batzuetan koma+juntagailua egitura
Enumerazioetan, osagaien artean koma	Enumerazioetan, osagaien artean koma
Deikiak komaz markatu	-
Estilo zuzeneko esaldiak komaz bereizi	Estilo zuzeneko esaldiak komaz bereizi
Aposizio ez-murritztaileak koma artean	Aposizio ez-murritztaileak koma artean
Tartekiak koma artean	Tartekiak koma artean
Lokailuak eta diskurtso-antolatzaileak koma artean	Izenen ondoko modifikatzaile ez-murritztaileak eta tartekiak koma artean
Aditz nagusiaren aurreko mintzagaiaren ondoren, koma	Esaldiaren hasieran, sarbide gisa doazen sintagmak edo perpausak komaz mugatuta
Aditz nagusiaren ondorengo perpaus zirkunstantzialen aurretik, koma	Esaldi bukaerako elementu osagarriak komaz bereizi
Aditza isildua dagoenean, komaz markatu	-

Taula: Koma jartzeko arauen konparazioa: euskara vs ingelesa

## Komen zuzenketa hizkuntzaren ezagutzan oinarrituta

- Komak zuzentzeko 19 CG erregela:

```
MAP (&OKER_KOMA_FALTA_1_1)  
TARGET EDOZEIN_KAT  
IF (1 BAINA + JNT);
```

## Komen zuzenketa hizkuntzaren ezagutzan oinarrituta

- Komak zuzentzeko 19 CG erregela:

MAP (&OKER\_KOMA\_FALTA\_1\_1)  
TARGET EDOZEIN\_KAT  
IF (1 BAINA + JNT);

- Adibidez:

*"Erabakia ez da bete beharrekoa baina ondorio politiko eta juridikoak izango ditu."*

## Komen zuzenketa hizkuntzaren ezagutzan oinarrituta

- Komak zuzentzeko 19 CG erregela:

MAP (&OKER\_KOMA\_FALTA\_1\_1)  
TARGET EDOZEIN\_KAT  
IF (1 BAINA + JNT);

- Adibidez:

*"Erabakia ez da bete beharrekoa, baina ondorio politiko eta juridikoak izango ditu."*

## Komen zuzenketa hizkuntzaren ezagutzan oinarrituta

- Komak zuzentzeko 19 CG erregela:

MAP (&OKER\_KOMA\_FALTA\_1\_1)  
TARGET EDOZEIN\_KAT  
IF (1 BAINA + JNT);

- Adibidez:

*"Erabakia ez da bete beharrekoa, baina ondorio politiko eta juridikoak izango ditu."*

- Emaitzak:

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>Erregeletan oinarrituta</b>	93,1	96,7	94,9	56,9	27,2	36,8

**Taula:** CG erregelekin lortutako komen berreskuratze-emaitzak. 0 klasea: atzetik komarik ez daramaten tokenak; 1 klasea: atzetik koma daramaten tokenak.



## Komen zuzenketa hizkuntzaren ezagutzan oinarrituta

- Komak zuzentzeko 19 CG erregela:

MAP (&OKER\_KOMA\_FALTA\_1\_1)  
TARGET EDOZEIN\_KAT  
IF (1 BAINA + JNT);

- Adibidez:

*"Erabakia ez da bete beharrekoa, baina ondorio politiko eta juridikoak izango ditu."*

- Emaitzak:

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>Erregeletan oinarrituta</b>	93,1	96,7	94,9	56,9	27,2	<b>36,8</b>

**Taula:** CG erregelekin lortutako komen berreskuratze-emaitzak. 0 klasea: atzetik komarik ez daramaten tokenak; 1 klasea: atzetik koma daramaten tokenak.

## Komen zuzenketa hizkuntzaren ezagutzan oinarrituta

- Komak zuzentzeko 19 CG erregela:

MAP (&OKER\_KOMA\_FALTA\_1\_1)  
TARGET EDOZEIN\_KAT  
IF (1 BAINA + JNT);

- Adibidez:

*"Erabakia ez da bete beharrekoa, baina ondorio politiko eta juridikoak izango ditu."*

- Emaitzak:

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>Erregeletan oinarrituta</b>	93,1	96,7	94,9	<b>56,9</b>	27,2	36,8

**Taula:** CG erregelekin lortutako komen berreskuratze-emaitzak. 0 klasea: atzetik komarik ez daramaten tokenak; 1 klasea: atzetik koma daramaten tokenak.

## Komen zuzenketa ikasketa automatikoan oinarrituta

- Esperimentuen prestaketa
- Egindako saioak
- Komen zuzenketa erregelak eta ikasketa automatikoa konbinatuz
- Jatorrizko komen eragina saihesten
- Ebaluazio kualitatiboa
- Erroreen analisia

## Esperimentuen prestaketa (I)

- Ikasketa-corpus gisa, komak ondo jarrita dituen edozein corpus erabil daiteke:

	<b>Ikasketarako</b>	<b>Garapenerako</b>	<b>Testerako</b>	<b>Guztira</b>
<b>Tokenak</b>	101.250	28.500	5.250	135.000

**Taula:** Komak ikasteko erabilitako *Euskaldunon Egunkariako* corpusa

## Esperimentuen prestaketa (I)

- Ikasketa-corpus gisa, komak ondo jarrita dituen edozein corpus erabil daiteke:

	<b>Ikasketarako</b>	<b>Garapenerako</b>	<b>Testerako</b>	<b>Guztira</b>
<b>Tokenak</b>	101.250	28.500	5.250	135.000

**Taula:** Komak ikasteko erabilitako *Euskaldunon Egunkariako* corpusa

- Oinarrizko neurriak:

	<b>0</b>			<b>1</b>		
	<b>Doit.</b>	<b>Est.</b>	$F_1$	<b>Doit.</b>	<b>Est.</b>	$F_1$
<b>baseline_%8</b>	92,7	92,4	92,6	7,6	7,9	7,8
<b>baseline_100</b>	94,1	80,2	86,6	12,5	35,8	18,5
<b>baseline_200</b>	94,4	75,6	84,0	12,1	42,7	18,9
<b>baseline_300</b>	94,5	72,4	82,0	11,2	46,5	18,7
<b>baseline_ikasketakoak</b>	94,6	55,6	70,0	9,6	59,6	16,5

## Esperimentuen prestaketa (I)

- Ikasketa-corpus gisa, komak ondo jarrita dituen edozein corpus erabil daiteke:

	<b>Ikasketarako</b>	<b>Garapenerako</b>	<b>Testerako</b>	<b>Guztira</b>
<b>Tokenak</b>	101.250	28.500	5.250	135.000

**Taula:** Komak ikasteko erabilitako *Euskaldunon Egunkariako* corpusa

- Oinarrizko neurriak:

	<b>0</b>			<b>1</b>		
	<b>Doit.</b>	<b>Est.</b>	$F_1$	<b>Doit.</b>	<b>Est.</b>	$F_1$
<b>baseline_%8</b>	92,7	92,4	<b>92,6</b>	7,6	7,9	7,8
<b>baseline_100</b>	94,1	80,2	86,6	12,5	35,8	18,5
<b>baseline_200</b>	94,4	75,6	84,0	12,1	42,7	18,9
<b>baseline_300</b>	94,5	72,4	82,0	11,2	46,5	18,7
<b>baseline_ikasketakoak</b>	94,6	55,6	70,0	9,6	59,6	16,5

## Esperimentuen prestaketa (I)

- Ikasketa-corpus gisa, komak ondo jarrita dituen edozein corpus erabil daiteke:

	Ikasketarako	Garapenerako	Testerako	Guztira
<b>Tokenak</b>	101.250	28.500	5.250	135.000

**Taula:** Komak ikasteko erabilitako *Euskaldunon Egunkariako* corpusa

- Oinarrizko neurriak:

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>baseline_%8</b>	92,7	92,4	<b>92,6</b>	7,6	7,9	7,8
<b>baseline_100</b>	94,1	80,2	86,6	12,5	35,8	18,5
<b>baseline_200</b>	94,4	75,6	84,0	12,1	42,7	<b>18,9</b>
<b>baseline_300</b>	94,5	72,4	82,0	11,2	46,5	18,7
<b>baseline_ikasketakoak</b>	94,6	55,6	70,0	9,6	59,6	16,5

## Esperimentuen prestaketa (II)

- Ikasi beharrekoa  $\Rightarrow$  token bakoitzaren ondoren koma datorren ala ez. Horretarako:
  - Corpusetik jatorrizko komak kendu
  - Token bakoitzari azken atributu gisa gehitu: atzetik koma daraman ala ez



## Esperimentuen prestaketa (III)

```
@RELATION komak.eu  
@ATTRIBUTE hitza REAL  
@ATTRIBUTE lema REAL  
@ATTRIBUTE kat {-,ADB,ADI,ADJ,...}  
...  
@ATTRIBUTE puntu {0,1}  
@ATTRIBUTE bi_puntu {0,1}  
...  
@ATTRIBUTE zenbat_ADK_ezk REAL  
@ATTRIBUTE zenbat_ADK_esk REAL  
...  
@ATTRIBUTE koma {0,1}
```

## Esperimentuen prestaketa (III)

```

@RELATION komak.eu
@ATTRIBUTE hitza REAL
@ATTRIBUTE lema REAL
@ATTRIBUTE kat {-,ADB,ADI,ADJ,...}
...
@ATTRIBUTE puntu {0,1}
@ATTRIBUTE bi_puntu {0,1}
...
@ATTRIBUTE zenbat_ADK_ek REAL
@ATTRIBUTE zenbat_ADK_esk REAL
...
@ATTRIBUTE koma {0,1}

@DATA
380,235,ADB,ADOARR,-,-,0,0,1,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1,0,4,0,1,1,18,0
2089,1054,LOT,LOK,-,-,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1,1,4,0,1,2,17,1
638,1412,IZE,ARR,-,-,0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1,1,3,0,1,3,16,0
6459,669,ADJ,IZO,ABS,-,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1,1,3,0,1,4,15,0
13,666,IZE,ARR,ERG,-,0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1,2,2,0,1,6,13,0
.....

```

## Egindako saioak: leihoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>(-2,+5)</b>	95,6	98,2	96,9	64,8	43,1	51,8
<b>(-3,+5)</b>	95,7	97,9	96,8	62,7	44,1	51,8
<b>(-4,+5)</b>	95,7	98,0	96,8	63,4	44,6	52,0
<b>(-5,+5)</b>	95,5	98,1	96,8	63,5	41,7	50,3
<b>(-5,+4)</b>	95,5	98,2	96,8	64,0	41,7	50,5
<b>(-5,+3)</b>	95,6	98,1	96,9	64,3	43,2	51,7
<b>(-5,+2)</b>	95,6	98,2	96,9	65,0	42,4	51,4
<b>(-6,+2)</b>	95,6	98,2	96,9	64,5	42,1	50,9
<b>(-6,+3)</b>	95,6	98,2	96,9	64,6	42,6	51,4
<b>(-8,+2)</b>	95,6	98,2	96,9	64,5	42,5	51,3
<b>(-8,+3)</b>	95,6	97,9	96,7	61,5	43,1	50,7
<b>(-8,+8)</b>	95,6	97,8	96,7	60,4	42,2	49,7

Taula: Leihoaren aukeraketa

## Egindako saioak: leihoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
(-2,+5)	95,6	98,2	<b>96,9</b>	64,8	43,1	51,8
(-3,+5)	95,7	97,9	<b>96,8</b>	62,7	44,1	51,8
(-4,+5)	95,7	98,0	<b>96,8</b>	63,4	44,6	52,0
(-5,+5)	95,5	98,1	<b>96,8</b>	63,5	41,7	50,3
(-5,+4)	95,5	98,2	<b>96,8</b>	64,0	41,7	50,5
(-5,+3)	95,6	98,1	<b>96,9</b>	64,3	43,2	51,7
(-5,+2)	95,6	98,2	<b>96,9</b>	65,0	42,4	51,4
(-6,+2)	95,6	98,2	<b>96,9</b>	64,5	42,1	50,9
(-6,+3)	95,6	98,2	<b>96,9</b>	64,6	42,6	51,4
(-8,+2)	95,6	98,2	<b>96,9</b>	64,5	42,5	51,3
(-8,+3)	95,6	97,9	<b>96,7</b>	61,5	43,1	50,7
(-8,+8)	95,6	97,8	<b>96,7</b>	60,4	42,2	49,7

Taula: Leihoaren aukeraketa

## Egindako saioak: leihoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>(-2,+5)</b>	95,6	98,2	96,9	64,8	43,1	51,8
<b>(-3,+5)</b>	95,7	97,9	96,8	62,7	44,1	51,8
<b>(-4,+5)</b>	95,7	98,0	96,8	63,4	44,6	<b>52,0</b>
<b>(-5,+5)</b>	95,5	98,1	96,8	63,5	41,7	50,3
<b>(-5,+4)</b>	95,5	98,2	96,8	64,0	41,7	50,5
<b>(-5,+3)</b>	95,6	98,1	96,9	64,3	43,2	51,7
<b>(-5,+2)</b>	95,6	98,2	96,9	65,0	42,4	51,4
<b>(-6,+2)</b>	95,6	98,2	96,9	64,5	42,1	50,9
<b>(-6,+3)</b>	95,6	98,2	96,9	64,6	42,6	51,4
<b>(-8,+2)</b>	95,6	98,2	96,9	64,5	42,5	51,3
<b>(-8,+3)</b>	95,6	97,9	96,7	61,5	43,1	50,7
<b>(-8,+8)</b>	95,6	97,8	96,7	60,4	42,2	49,7

Taula: Leihoaren aukeraketa

## Egindako saioak: leihoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>(-2,+5)</b>	95,6	98,2	96,9	64,8	43,1	51,8
<b>(-3,+5)</b>	95,7	97,9	96,8	62,7	44,1	51,8
<b>(-4,+5)</b>	95,7	98,0	96,8	63,4	44,6	<b>52,0</b>
<b>(-5,+5)</b>	95,5	98,1	96,8	63,5	41,7	50,3
<b>(-5,+4)</b>	95,5	98,2	96,8	64,0	41,7	50,5
<b>(-5,+3)</b>	95,6	98,1	96,9	64,3	43,2	51,7
<b>(-5,+2)</b>	95,6	98,2	96,9	<b>65,0</b>	42,4	51,4
<b>(-6,+2)</b>	95,6	98,2	96,9	64,5	42,1	50,9
<b>(-6,+3)</b>	95,6	98,2	96,9	64,6	42,6	51,4
<b>(-8,+2)</b>	95,6	98,2	96,9	64,5	42,5	51,3
<b>(-8,+3)</b>	95,6	97,9	96,7	61,5	43,1	50,7
<b>(-8,+8)</b>	95,6	97,8	96,7	60,4	42,2	49,7

Taula: Leihoaren aukeraketa

## Egindako saioak: leihoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>(-2,+5)</b>	95,6	98,2	96,9	64,8	43,1	51,8
<b>(-3,+5)</b>	95,7	97,9	96,8	62,7	44,1	51,8
<b>(-4,+5)</b>	95,7	98,0	96,8	<b>63,4</b>	44,6	<b>52,0</b>
<b>(-5,+5)</b>	95,5	98,1	96,8	63,5	41,7	50,3
<b>(-5,+4)</b>	95,5	98,2	96,8	64,0	41,7	50,5
<b>(-5,+3)</b>	95,6	98,1	96,9	64,3	43,2	51,7
<b>(-5,+2)</b>	95,6	98,2	96,9	<b>65,0</b>	42,4	<b>51,4</b>
<b>(-6,+2)</b>	95,6	98,2	96,9	64,5	42,1	50,9
<b>(-6,+3)</b>	95,6	98,2	96,9	64,6	42,6	51,4
<b>(-8,+2)</b>	95,6	98,2	96,9	64,5	42,5	51,3
<b>(-8,+3)</b>	95,6	97,9	96,7	61,5	43,1	50,7
<b>(-8,+8)</b>	95,6	97,8	96,7	60,4	42,2	49,7

Taula: Leihoaren aukeraketa

## Egindako saioak: algoritmoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>C4.5</b>	95,6	98,2	96,9	65,2	42,4	51,4
<b>Naive Bayes</b>	94,8	95,6	95,2	37,6	33,5	35,5
<b>SVM</b>	93,6	99,4	96,5	67,2	14,3	23,6

Taula: Ikasketa-algoritmoaren aukeraketa



## Egindako saioak: algoritmoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>C4.5</b>	95,6	98,2	96,9	65,2	42,4	<b>51,4</b>
<b>Naive Bayes</b>	94,8	95,6	95,2	37,6	33,5	35,5
<b>SVM</b>	93,6	99,4	96,5	67,2	14,3	23,6

Taula: Ikasketa-algoritmoaren aukeraketa

## Egindako saioak: algoritmoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>C4.5</b>	95,6	98,2	96,9	65,2	42,4	<b>51,4</b>
<b>Naive Bayes</b>	94,8	95,6	95,2	37,6	33,5	35,5
<b>SVM</b>	93,6	99,4	96,5	<b>67,2</b>	14,3	23,6

Taula: Ikasketa-algoritmoaren aukeraketa

## Egindako saioak: algoritmoaren aukeraketa

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>C4.5</b>	95,6	98,2	96,9	<b>65,2</b>	42,4	<b>51,4</b>
<b>Naive Bayes</b>	94,8	95,6	95,2	37,6	33,5	35,5
<b>SVM</b>	93,6	99,4	96,5	<b>67,2</b>	14,3	<b>23,6</b>

Taula: Ikasketa-algoritmoaren aukeraketa

## Egindako saioak: corpus motaren eragina

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>Egunkaria</b>	93,0	97,9	95,4	48,7	21,4	29,7
<b>Filosofia-itzulpena</b>	93,0	97,0	95,0	60,4	38,8	47,3
<b>Literatura</b>	92,7	97,5	95,0	50,4	25,0	33,4
<b>Zientzia eta teknika</b>	94,9	98,5	96,6	49,5	22,4	30,8

**Taula:** *Cross-validation* emaitzak corpus motaren arabera (25.000 tokeneko corpusak)

## Egindako saioak: corpus motaren eragina

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>Egunkaria</b>	93,0	97,9	95,4	48,7	21,4	29,7
<b>Filosofia-itzulpena</b>	93,0	97,0	95,0	60,4	38,8	<b>47,3</b>
<b>Literatura</b>	92,7	97,5	95,0	50,4	25,0	33,4
<b>Zientzia eta teknika</b>	94,9	98,5	96,6	49,5	22,4	30,8

**Taula:** *Cross-validation* emaitzak corpus motaren arabera (25.000 tokeneko corpusak)

## Egindako saioak: ingeleseko corpusarekin

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
Euskara	95,6	98,2	96,9	65,2	42,4	51,4
Ingelesa	97,8	99,7	98,7	83,3	38,7	52,8

Taula: Hizkuntzaren araberrako emaitzak: euskara vs ingelesa

## Egindako saioak: ingeleseko corpusarekin

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>Euskara</b>	95,6	98,2	<b>96,9</b>	65,2	42,4	<b>51,4</b>
<b>Ingelesa</b>	97,8	99,7	<b>98,7</b>	83,3	38,7	<b>52,8</b>

Taula: Hizkuntzaren araberrako emaitzak: euskara vs ingelesa

## Egindako saioak: atributu berriekin

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>Atributu berririk gabe</b>	95,6	98,2	96,9	65,2	42,4	51,4
<b>(1) 300 atributu berri</b>	96,0	98,3	97,2	69,6	48,6	57,2
<b>(2) 3 atributu berri</b>	96,0	98,1	97,0	66,5	48,1	55,8

**Taula:** Komen aurretik maizen agertzen diren hitzak atributu gisa gehituz

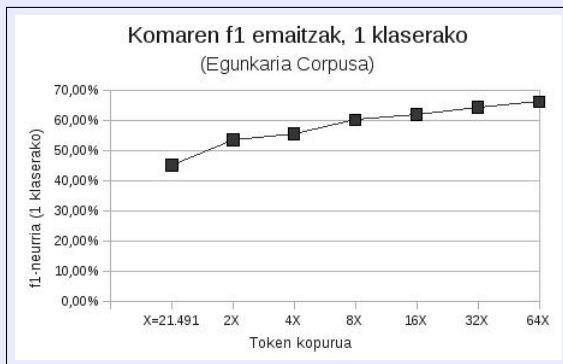


## Egindako saioak: atributu berriekin

	0			1		
	Doit.	Est.	$F_1$	Doit.	Est.	$F_1$
<b>Atributu berririk gabe</b>	95,6	98,2	96,9	65,2	42,4	51,4
<b>(1) 300 atributu berri</b>	96,0	98,3	97,2	<b>69,6</b>	<b>48,6</b>	<b>57,2</b>
<b>(2) 3 atributu berri</b>	96,0	98,1	97,0	66,5	48,1	55,8

**Taula:** Komen aurretik maizen agertzen diren hitzak atributu gisa gehituz

## Egindako saioak: corpusaren tamainaren eragina



**Irudia:** Corpusaren tamainaren eragina (hamar zatiko *cross-validation* baliatuta)

## Egindako saioak: kateen eta perpausen informazioa gehituz

<b>Ikasketa-corpora</b>	<b>Test-corpora</b>	<b>Desanb.</b>	$F_1$ neurria
<i>Komaduna</i>	<i>Komaduna</i>	autom.	83,17
<i>Komagabea</i>	<i>Komagabea</i>	autom.	82,24

Taula: Komaren eragina *FR-P* bidezko **kate-identifikatzailean**

## Egindako saioak: kateen eta perpausen informazioa gehituz

Ikasketa-corpora	Test-corpora	Desanb.	$F_1$ neurria
<i>Komaduna</i>	<i>Komaduna</i>	autom.	83,17
<i>Komagabea</i>	<i>Komagabea</i>	autom.	82,24

Taula: Komaren eragina *FR-P* bidezko **kate-identifikatzailean**

Ikasketa-corpora	Test-corpora	Desanb.	$F_1$ neurria
<i>Komaduna</i>	<i>Komaduna</i>	autom.	77,24
<i>Komagabea</i>	<i>Komagabea</i>	autom.	73,66

Taula: Komaren eragina *FR-P* bidezko **perpaus-identifikatzailean**

## Egindako saioak: kateen eta perpausen informazioa gehituz

	1		
	Doit.	Est.	$F_1$
<b>Kate-info. eta perpaus-info. gabe</b>	69,6	48,6	57,2
<b>Kate-ident komagabearen info. gehituta</b>	70,4	48,5	57,4
<b>Kate- eta perpaus-ident komagabeen info. gehituta</b>	76,6	55,7	64,5
<b>Kate-ident komadunaren info. gehituta</b>	73,0	50,7	59,8
<b>Kate- eta perpaus-ident komadunen info. gehituta</b>	78,4	59,8	67,9

**Taula:** Komen zuzenketa, kate- eta perpaus-identifikatzaile komadunen edo komagabeen informazioa gehituta

## Egindako saioak: kateen eta perpausen informazioa gehituz

	1		
	Doit.	Est.	$F_1$
<b>Kate-info. eta perpaus-info. gabe</b>	69,6	48,6	57,2
<b>Kate-ident komagabearen info. gehituta</b>	70,4	48,5	<b>57,4</b>
<b>Kate- eta perpaus-ident komagabeen info. gehituta</b>	76,6	55,7	64,5
<b>Kate-ident komadunaren info. gehituta</b>	73,0	50,7	59,8
<b>Kate- eta perpaus-ident komadunen info. gehituta</b>	78,4	59,8	67,9

**Taula:** Komen zuzenketa, kate- eta perpaus-identifikatzaile komadunen edo komagabeen informazioa gehituta

## Egindako saioak: kateen eta perpausen informazioa gehituz

	1		
	Doit.	Est.	$F_1$
Kate-info. eta perpaus-info. gabe	69,6	48,6	57,2
Kate-ident komagabearen info. gehituta	70,4	48,5	57,4
Kate- eta perpaus-ident komagabeen info. gehituta	76,6	55,7	<b>64,5</b>
Kate-ident komadunaren info. gehituta	73,0	50,7	59,8
Kate- eta perpaus-ident komadunen info. gehituta	78,4	59,8	67,9

**Taula:** Komen zuzenketa, kate- eta perpaus-identifikatzaile komadunen edo komagabeen informazioa gehituta

## Egindako saioak: kateen eta perpausen informazioa gehituz

	1		
	Doit.	Est.	$F_1$
Kate-info. eta perpaus-info. gabe	69,6	48,6	57,2
Kate-ident komagabearen info. gehituta	70,4	48,5	57,4
Kate- eta perpaus-ident komagabeen info. gehituta	76,6	55,7	64,5
Kate-ident komadunaren info. gehituta	73,0	50,7	<b>59,8</b>
Kate- eta perpaus-ident komadunen info. gehituta	78,4	59,8	67,9

**Taula:** Komen zuzenketa, kate- eta perpaus-identifikatzaile komadunen edo komagabeen informazioa gehituta



## Egindako saioak: kateen eta perpausen informazioa gehituz

	1		
	Doit.	Est.	$F_1$
<b>Kate-info. eta perpaus-info. gabe</b>	69,6	48,6	57,2
<b>Kate-ident komagabearen info. gehituta</b>	70,4	48,5	57,4
<b>Kate- eta perpaus-ident komagabeen info. gehituta</b>	76,6	55,7	64,5
<b>Kate-ident komadunaren info. gehituta</b>	73,0	50,7	59,8
<b>Kate- eta perpaus-ident komadunen info. gehituta</b>	78,4	59,8	<b>67,9</b>

**Taula:** Komen zuzenketa, kate- eta perpaus-identifikatzaile komadunen edo komagabeen informazioa gehituta

## Komen zuzenketa erregelak eta ikasketa automatikoa konbinatuz

	1		
	Doit.	Est.	$F_1$
<b>CG erregelak</b>	56,9	27,2	36,8
<b>Ikasketa automatikoa KPI-komagabearekin</b>	76,6	55,7	64,5
<b>CG erregelak + ikask. autom. KPI-komagabearekin</b>	77,8	55,0	64,4
<b>Ikask. autom. KPI-komadunarekin</b>	78,4	59,8	67,9
<b>CG erregelak + ikask. autom. KPI-komadunarekin</b>	79,0	61,4	69,1

Taula: Erregelaren informazioa gehitzea, pilaratzearen bidez

## Komen zuzenketa erregelak eta ikasketa automatikoa konbinatuz

	1		
	Doit.	Est.	$F_1$
CG erregelak	56,9	27,2	36,8
<b>Ikasketa automatikoa KPI-komagabearekin</b>	<b>76,6</b>	55,7	<b>64,5</b>
<b>CG erregelak + ikask. autom. KPI-komagabearekin</b>	<b>77,8</b>	55,0	<b>64,4</b>
Ikask. autom. KPI-komadunarekin	78,4	59,8	67,9
<b>CG erregelak + ikask. autom. KPI-komadunarekin</b>	79,0	61,4	69,1

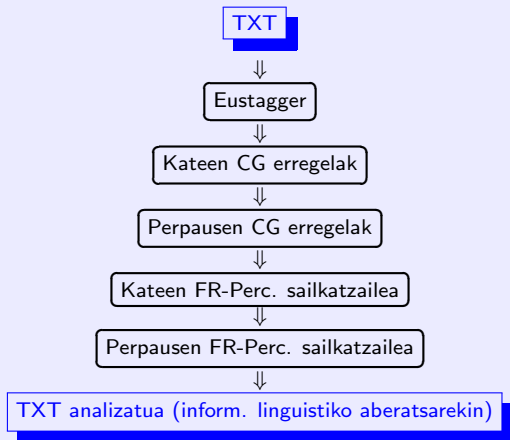
Taula: Erregelaren informazioa gehitzea, pilaratzearen bidez

## Komen zuzenketa erregelak eta ikasketa automatikoa konbinatuz

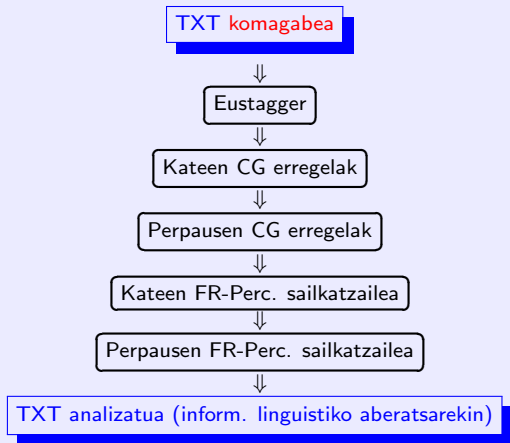
	1		
	Doit.	Est.	$F_1$
CG erregelak	56,9	27,2	36,8
Ikasketa automatikoa KPI-komagabearekin	76,6	55,7	64,5
CG erregelak + ikask. autom. KPI-komagabearekin	77,8	55,0	64,4
Ikask. autom. KPI-komadunarekin	78,4	59,8	<b>67,9</b>
CG erregelak + ikask. autom. KPI-komadunarekin	79,0	61,4	<b>69,1</b>

Taula: Erregelaren informazioa gehitzea, pilaratzearen bidez

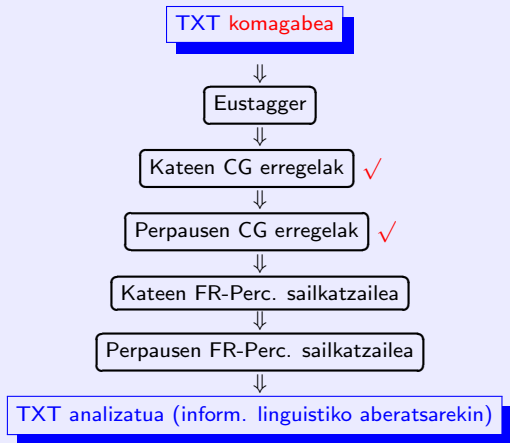
## Jatorrizko komen eragina saihesten



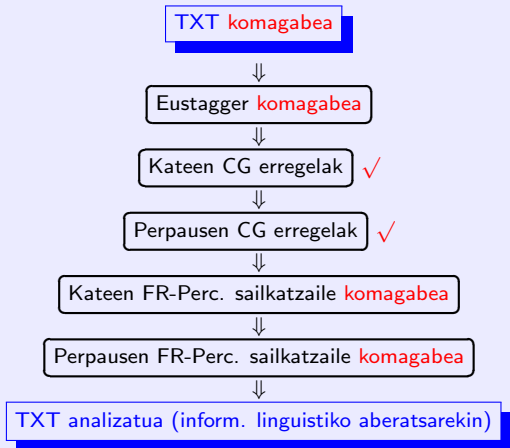
## Jatorrizko komen eragina saihesten



# Jatorrizko komen eragina saihesten

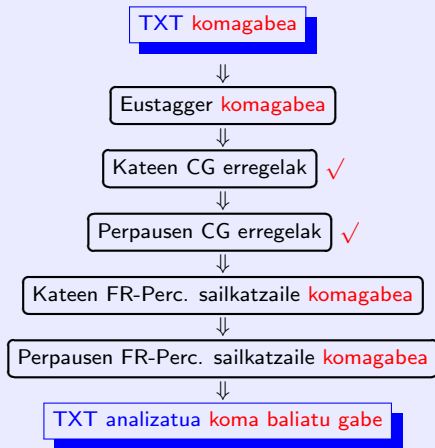


## Jatorrizko komen eragina saihesten





# Jatorrizko komen eragina saihesten



## Jatorrizko komen eragina saihesten

	1		
	Doit.	Est.	$F_1$
<b>Corpus komaduna + Eustagger komaduna</b>	77,8	55,0	64,4
<b>Corpus komagabea + Eustagger komagabea</b>	69,3	33,3	45,0

**Taula:** Eustagger komagabeak edo komadunak emandako informazio linguistikoaren arabera

## Jatorrizko komen eragina saihesten

	1		
	Doit.	Est.	$F_1$
Corpus komaduna + Eustagger komaduna	77,8	55,0	<b>64,4</b>
Corpus komagabea + Eustagger komagabea	69,3	33,3	<b>45,0</b>

**Taula:** Eustagger komagabeak edo komadunak emandako informazio linguistikoaren arabera

## Jatorrizko komen eragina saihesten

	1		
	Doit.	Est.	$F_1$
Corpus komaduna + Eustagger komaduna	77,8	55,0	<b>64,4</b>
Corpus komagabea + Eustagger komagabea	<b>69,3</b>	33,3	<b>45,0</b>

**Taula:** Eustagger komagabeak edo komadunak emandako informazio linguistikoaren arabera

## Ebaluazio kualitatiboa

	1		
	Doit.	Est.	$F_1$
<b>Ikask. autom. KPI-komagabe + CG erregelak</b>	77,6	52,7	62,8
<b>Hizkuntzalari1</b>	79,1	85,9	82,3
<b>Hizkuntzalari2</b>	76,1	76,4	76,3

**Taula:** Test-corpuseko emaitzak vs hizkuntzalarien etiketatzearen emaitzak

## Ebaluazio kualitatiboa

	1		
	Doit.	Est.	$F_1$
<b>Ikask. autom. KPI-komagabe + CG erregelak</b>	77,6	52,7	<b>62,8</b>
<b>Hizkuntzalari1</b>	79,1	85,9	82,3
<b>Hizkuntzalari2</b>	76,1	76,4	<b>76,3</b>

**Taula:** Test-corpuseko emaitzak vs hizkuntzalarien etiketatzearen emaitzak

## Ebaluazio kualitatiboa

	1		
	Doit.	Est.	$F_1$
<b>Ikask. autom. KPI-komagabe + CG erregelak</b>	<b>77,6</b>	52,7	62,8
<b>Hizkuntzalari1</b>	79,1	85,9	82,3
<b>Hizkuntzalari2</b>	<b>76,1</b>	76,4	76,3

**Taula:** Test-corpuseko emaitzak vs hizkuntzalarien etiketatzearen emaitzak

## Ebaluazio kualitatiboa

	1		
	Doit.	Est.	$F_1$
Ikask. autom. KPI-komagabe + CG erregelak	77,6	52,7	62,8
Hizkuntzalari1	79,1	85,9	82,3
Hizkuntzalari2	76,1	76,4	<b>76,3</b>

Taula: Test-corpuseko emaitzak vs hizkuntzalarien etiketatzearen emaitzak

	1		
	Doit.	Est.	$F_1$
Hizkuntzalari1, hizkuntzalari2-rekiko	73,8	79,7	<b>76,6</b>

Taula: Hizkuntzalarien arteko adostasuna



## Ebaluazio kualitatiboa

	1		
	Doit.	Est.	$F_1$
Ikask. autom. KPI-komagabearekin + CG erregelak	77,6	52,7	62,8
<b>Ebaluazio kualitatiboa (hiru erref.)</b>	<b>83,01</b>	58,46	<b>68,61</b>

**Taula:** Sailkatzailearen bateragarritasuna test-corpusarekiko eta hiru erreferentziekiko

## Erroren analisia

- 1 Gurean igandeko egunkariak aste osoa ematen dute komunikatibitate **&SOBRAN** bueltaka eta jiraka **&FALTAN** orain toalleroan **&FALTAN** orain erradiadorean.
- 2 Pasa ziren egiazko ospakizunak eta itxurazkoak **&FALTAN** etorri ziren lehen adierazpenak eta hasierako azterketak **&KOMA** argazkiak eta gezur ezkutatuak.
- 3 Batzuen ustez **&KOMA** inora ez doana **&KOMA** ezer egiten ez duena **&KOMA** eta beste batzuentzat **&FALTAN** gehiegi egiten duena **&KOMA** urrutiegi eta arinegi doana.
- 4 UEUk **&FALTAN** EIREk **&FALTAN** Euskal Adarrak **&KOMA** Barrutiak eta Uniekimenak "eztabaidaren erdigunean" jarri nahi dute aldarrikapena.

# Aurkezpenaren eskema

- 1 Testuingurua
- 2 Kateen eta perpausen identifikazioa
- 3 Komaren zuzenketa automatikoa
- 4 Ondorioak eta etorkizuneko lanak
  - Ekarpenak
  - Ondorioak
  - Etorkizuneko lanak
  - Argitalpenak

## Ekarpen garrantzitsuenak

- Euskarako kate- eta perpaus-identifikatzailearen garapena

## Ekarpen garrantzitsuenak

- Euskarako kate- eta perpaus-identifikatzailearen garapena
  - Analisi sintaktiko osoa bideratzeko
  - HPko hainbat arlotarako: informazioaren erauzketa, laburpenen sorkuntza, itzulpen automatikoa...
  - Koma-zuzentzaileerako

## Ekarpen garrantzitsuenak

- Euskarako kate- eta perpaus-identifikatzailearen garapena
  - Analisi sintaktiko osoa bideratzeko
  - HPko hainbat arlotarako: informazioaren erauzketa, laburpenen sorkuntza, itzulpen automatikoa...
  - Koma-zuzentzaileerako
- Euskarako koma-zuzentzaile automatikoaren garapena

## Ekarpen garrantzitsuenak

- Euskarako kate- eta perpaus-identifikatzailearen garapena
  - Analisi sintaktiko osoa bideratzeko
  - HPko hainbat arlotarako: informazioaren erauzketa, laburpenen sorkuntza, itzulpen automatikoa. . .
  - Koma-zuzentzaileako
- Euskarako koma-zuzentzaile automatikoaren garapena
  - Estilo- eta gramatika-zuzentzailean txertatzeko
  - Analizatzaile eta desanbiguatzaile sintaktiko automatikoak hobetzeko
  - Ahotsaren ezagutza-sistemetan integratzeko

## Kate-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
<b>Erreg</b>	-	Autom.	51,48
<b>Baseline-a</b>	-	Autom.	52,00
<b>FR-P oin</b>	% 25	Autom.	69,58
<b>FR-P oin + dek</b>	% 25	Autom.	78,77
<b>FR-P oin + dek + Erreg</b>	% 25	Autom.	79,21
<b>FR-P oin + dek + Erreg</b>	% 100	Autom.	82,64
<b>FR-P oin + dek + Erreg</b>	% 100	Eskuz	90,52

Taula: Garapen-corpusean egindako ebaluazioa



## Kate-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
Erreg	-	Autom.	<b>51,48</b>
Baseline-a	-	Autom.	<b>52,00</b>
FR-P oin	% 25	Autom.	<b>69,58</b>
FR-P oin + dek	% 25	Autom.	78,77
FR-P oin + dek + Erreg	% 25	Autom.	79,21
FR-P oin + dek + Erreg	% 100	Autom.	82,64
FR-P oin + dek + Erreg	% 100	Eskuz	90,52

Taula: Garapen-corpusean egindako ebaluazioa

## Kate-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
Erreg	-	Autom.	51,48
Baseline-a	-	Autom.	52,00
FR-P oin	% 25	Autom.	<b>69,58</b>
FR-P oin + dek	% 25	Autom.	<b>78,77</b>
FR-P oin + dek + Erreg	% 25	Autom.	79,21
FR-P oin + dek + Erreg	% 100	Autom.	82,64
FR-P oin + dek + Erreg	% 100	Eskuz	90,52

Taula: Garapen-corpusean egindako ebaluazioa

## Kate-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
<b>Erreg</b>	-	Autom.	51,48
<b>Baseline-a</b>	-	Autom.	52,00
<b>FR-P oin</b>	% 25	Autom.	69,58
<b>FR-P oin + dek</b>	% 25	Autom.	<b>78,77</b>
<b>FR-P oin + dek + Erreg</b>	% 25	Autom.	<b>79,21</b>
<b>FR-P oin + dek + Erreg</b>	% 100	Autom.	82,64
<b>FR-P oin + dek + Erreg</b>	% 100	Eskuz	90,52

Taula: Garapen-corpusean egindako ebaluazioa

## Kate-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
Erreg	-	Autom.	51,48
Baseline-a	-	Autom.	52,00
FR-P oin	% 25	Autom.	69,58
FR-P oin + dek	% 25	Autom.	78,77
FR-P oin + dek + Erreg	% 25	Autom.	<b>79,21</b>
FR-P oin + dek + Erreg	% 100	Autom.	<b>82,64</b>
FR-P oin + dek + Erreg	% 100	Eskuz	90,52

Taula: Garapen-corpusean egindako ebaluazioa

## Kate-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
Erreg	-	Autom.	51,48
Baseline-a	-	Autom.	52,00
FR-P oin	% 25	Autom.	69,58
FR-P oin + dek	% 25	Autom.	78,77
FR-P oin + dek + Erreg	% 25	Autom.	79,21
FR-P oin + dek + Erreg	% 100	Autom.	<b>82,64</b>
FR-P oin + dek + Erreg	% 100	Eskuz	<b>90,52</b>

Taula: Garapen-corpusean egindako ebaluazioa

## Kate-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
<b>Erreg</b>	-	Autom.	51,48
<b>Baseline-a</b>	-	Autom.	52,00
<b>FR-P oin</b>	% 25	Autom.	69,58
<b>FR-P oin + dek</b>	% 25	Autom.	78,77
<b>FR-P oin + dek + Erreg</b>	% 25	Autom.	79,21
<b>FR-P oin + dek + Erreg</b>	% 100	Autom.	<b>82,64</b>
<b>FR-P oin + dek + Erreg</b>	% 100	Eskuz	90,52

Taula: Garapen-corpusean egindako ebaluazioa

Hizkuntza	Teknika	Desanbiguatua	$F_1$
<i>Euskara</i>	<i>FR-P oin + dek + Erreg</i>	Autom.	<b>83,17</b>
<i>Ingelesa</i>	<i>FR-P</i>	Autom.	93,74

Taula: Test-corpusean egindako ebaluazioa

## Perpaus-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
Erreg	-	Autom.	49,71
Baseline-a	-	Autom.	48,79
FR-P oin	% 25	Autom.	69,99
FR-P oin+ak+dek+l+m	% 25	Autom.	73,22
FR-P oin+ak+dek+l+m+Erreg	% 25	Autom.	74,54
FR-P oin+ak+dek+l+m+Erreg	% 100	Autom.	78,11
FR-P oin+ak+dek+l+m+Erreg	% 100	Eskuz	78,39

Taula: Garapen-corpusean egindako ebaluazioa

## Perpaus-identifikatzailea

	Ik. corpusa	Desanbiguatua	$F_1$
Erreg	-	Autom.	<b>49,71</b>
Baseline-a	-	Autom.	<b>48,79</b>
FR-P oin	% 25	Autom.	<b>69,99</b>
FR-P oin+ak+dek+l+m	% 25	Autom.	73,22
FR-P oin+ak+dek+l+m+Erreg	% 25	Autom.	74,54
FR-P oin+ak+dek+l+m+Erreg	% 100	Autom.	78,11
FR-P oin+ak+dek+l+m+Erreg	% 100	Eskuz	78,39

Taula: Garapen-corpusean egindako ebaluazioa



## Perpaus-identifikatzailea

	Ik. corpusa	Desanbiguatua	$F_1$
Erreg	-	Autom.	49,71
Baseline-a	-	Autom.	48,79
FR-P oin	% 25	Autom.	<b>69,99</b>
FR-P oin+ak+dek+l+m	% 25	Autom.	<b>73,22</b>
FR-P oin+ak+dek+l+m+Erreg	% 25	Autom.	74,54
FR-P oin+ak+dek+l+m+Erreg	% 100	Autom.	78,11
FR-P oin+ak+dek+l+m+Erreg	% 100	Eskuz	78,39

Taula: Garapen-corpusean egindako ebaluazioa

## Perpaus-identifikatzailea

	Ik. corpusa	Desanbiguatua	$F_1$
Erreg	-	Autom.	49,71
Baseline-a	-	Autom.	48,79
FR-P oin	% 25	Autom.	69,99
FR-P oin+ak+dek+l+m	% 25	Autom.	<b>73,22</b>
FR-P oin+ak+dek+l+m+Erreg	% 25	Autom.	<b>74,54</b>
FR-P oin+ak+dek+l+m+Erreg	% 100	Autom.	78,11
FR-P oin+ak+dek+l+m+Erreg	% 100	Eskuz	78,39

Taula: Garapen-corpusean egindako ebaluazioa

## Perpaus-identifikatzailea

	Ik. corpusa	Desanbiguatua	$F_1$
Erreg	-	Autom.	49,71
Baseline-a	-	Autom.	48,79
FR-P oin	% 25	Autom.	69,99
FR-P oin+ak+dek+l+m	% 25	Autom.	73,22
FR-P oin+ak+dek+l+m+Erreg	% 25	Autom.	<b>74,54</b>
FR-P oin+ak+dek+l+m+Erreg	% 100	Autom.	<b>78,11</b>
FR-P oin+ak+dek+l+m+Erreg	% 100	Eskuz	78,39

Taula: Garapen-corpusean egindako ebaluazioa

## Perpaus-identifikatzailea

	Ik. corpusa	Desanbiguatua	$F_1$
Erreg	-	Autom.	49,71
Baseline-a	-	Autom.	48,79
FR-P oin	% 25	Autom.	69,99
FR-P oin+ak+dek+l+m	% 25	Autom.	73,22
FR-P oin+ak+dek+l+m+Erreg	% 25	Autom.	74,54
FR-P oin+ak+dek+l+m+Erreg	% 100	Autom.	<b>78,11</b>
FR-P oin+ak+dek+l+m+Erreg	% 100	Eskuz	<b>78,39</b>

Taula: Garapen-corpusean egindako ebaluazioa

## Perpaus-identifikatzailea

	Ik. corpora	Desanbiguatua	$F_1$
Erreg	-	Autom.	49,71
Baseline-a	-	Autom.	48,79
FR-P oin	% 25	Autom.	69,99
FR-P oin+ak+dek+l+m	% 25	Autom.	73,22
FR-P oin+ak+dek+l+m+Erreg	% 25	Autom.	74,54
FR-P oin+ak+dek+l+m+Erreg	% 100	Autom.	<b>78,11</b>
FR-P oin+ak+dek+l+m+Erreg	% 100	Eskuz	78,39

Taula: Garapen-corpusean egindako ebaluazioa

Hizkuntza	Teknika	Desanbiguatua	$F_1$
<i>Euskara</i>	FR-P oin+ak+dek+l+m+e+Erreg	Autom.	<b>77,24</b>
<i>Ingelesa</i>	FR-P oin	Autom.	84,36

Taula: Test-corpusean egindako ebaluazioa

## Euskarako koma-zuzentzailea

- Komari buruzko azterketa teorikoa
- Euskara vs ingelesa
- Kate- eta perpaus-identifikatzaileak baliatu ditugu, koma-zuzentzailea hobetzeko
- Ebaluazio kualitatiboa egin dugu: 3 erreferentzia
- Eustagger-en koma erabiltzearen eragina aztertu dugu

## Euskarako koma-zuzentzailea

	1		
	Doit.	Est.	$F_1$
Baseline-200	12,1	42,7	18,9
<b>Err</b>	56,9	27,2	36,8
C4.5, -5+2	65,2	42,4	51,4
C4.5, -5+2, 300 atrib.	69,6	48,6	57,2
C4.5, -5+2, 300 atrib., KPI komagabearekin	76,6	55,7	64,5
Err + C4.5, -5+2, 300 atrib., KPI komagabearekin	77,8	55,0	64,4
C4.5, -5+2, 300 atrib., KPI komadunarekin	78,4	59,8	67,9
Err + C4.5, -5+2, 300 atrib., KPI komadunarekin	79,0	61,4	69,1
Err + C4.5, -5+2, 300 atrib., KPI eta Eustagger komagabeekin	69,3	33,3	45,0

Taula: Koma-zuzentzailearen emaitzen laburpena **garapen-corpusean**

## Euskarako koma-zuzentzailea

	1		
	Doit.	Est.	$F_1$
Baseline-200	12,1	42,7	<b>18,9</b>
<b>Err</b>	56,9	27,2	<b>36,8</b>
C4.5, -5+2	65,2	42,4	<b>51,4</b>
C4.5, -5+2, 300 atrib.	69,6	48,6	57,2
C4.5, -5+2, 300 atrib., KPI komagabearekin	76,6	55,7	64,5
Err + C4.5, -5+2, 300 atrib., KPI komagabearekin	77,8	55,0	64,4
C4.5, -5+2, 300 atrib., KPI komadunarekin	78,4	59,8	67,9
Err + C4.5, -5+2, 300 atrib., KPI komadunarekin	79,0	61,4	69,1
Err + C4.5, -5+2, 300 atrib., KPI eta Eustagger komagabeekin	69,3	33,3	45,0

Taula: Koma-zuzentzailearen emaitzen laburpena **garapen-corpusean**



## Euskarako koma-zuzentzailea

	1		
	Doit.	Est.	$F_1$
Baseline-200	12,1	42,7	18,9
<b>Err</b>	56,9	27,2	36,8
C4.5, -5+2	65,2	42,4	<b>51,4</b>
C4.5, -5+2, 300 atrib.	69,6	48,6	<b>57,2</b>
C4.5, -5+2, 300 atrib., KPI komagabearekin	76,6	55,7	64,5
Err + C4.5, -5+2, 300 atrib., KPI komagabearekin	77,8	55,0	64,4
C4.5, -5+2, 300 atrib., KPI komadunarekin	78,4	59,8	67,9
Err + C4.5, -5+2, 300 atrib., KPI komadunarekin	79,0	61,4	69,1
Err + C4.5, -5+2, 300 atrib., KPI eta Eustagger komagabeekin	69,3	33,3	45,0

Taula: Koma-zuzentzailearen emaitzen laburpena **garapen-corpusean**

## Euskarako koma-zuzentzailea

	1		
	Doit.	Est.	$F_1$
Baseline-200	12,1	42,7	18,9
Err	56,9	27,2	36,8
C4.5, -5+2	65,2	42,4	51,4
C4.5, -5+2, 300 atrib.	69,6	48,6	<b>57,2</b>
C4.5, -5+2, 300 atrib., KPI komagabearekin	76,6	55,7	<b>64,5</b>
Err + C4.5, -5+2, 300 atrib., KPI komagabearekin	<b>77,8</b>	55,0	64,4
C4.5, -5+2, 300 atrib., KPI komadunarekin	78,4	59,8	67,9
Err + C4.5, -5+2, 300 atrib., KPI komadunarekin	79,0	61,4	69,1
Err + C4.5, -5+2, 300 atrib., KPI eta Eustagger komagabeekin	69,3	33,3	45,0

Taula: Koma-zuzentzailearen emaitzen laburpena **garapen-corpusean**

## Euskarako koma-zuzentzailea

	1		
	Doit.	Est.	$F_1$
Baseline-200	12,1	42,7	18,9
<b>Err</b>	56,9	27,2	36,8
C4.5, -5+2	65,2	42,4	51,4
C4.5, -5+2, 300 atrib.	69,6	48,6	<b>57,2</b>
C4.5, -5+2, 300 atrib., KPI komagabearekin	76,6	55,7	64,5
Err + C4.5, -5+2, 300 atrib., KPI komagabearekin	77,8	55,0	64,4
C4.5, -5+2, 300 atrib., KPI komadunarekin	78,4	59,8	<b>67,9</b>
Err + C4.5, -5+2, 300 atrib., KPI komadunarekin	79,0	61,4	<b>69,1</b>
Err + C4.5, -5+2, 300 atrib., KPI eta Eustagger komagabeekin	69,3	33,3	45,0

Taula: Koma-zuzentzailearen emaitzen laburpena **garapen-corpusean**

## Euskarako koma-zuzentzailea

	1		
	Doit.	Est.	$F_1$
Baseline-200	12,1	42,7	18,9
Err	56,9	27,2	36,8
C4.5, -5+2	65,2	42,4	51,4
C4.5, -5+2, 300 atrib.	69,6	48,6	57,2
C4.5, -5+2, 300 atrib., KPI komagabearekin	76,6	55,7	64,5
Err + C4.5, -5+2, 300 atrib., KPI komagabearekin	77,8	55,0	<b>64,4</b>
C4.5, -5+2, 300 atrib., KPI komadunarekin	78,4	59,8	67,9
Err + C4.5, -5+2, 300 atrib., KPI komadunarekin	79,0	61,4	<b>69,1</b>
Err + C4.5, -5+2, 300 atrib., KPI eta Eustagger komagabeekin	69,3	33,3	<b>45,0</b>

Taula: Koma-zuzentzailearen emaitzen laburpena **garapen-corpusean**

## Euskarako koma-zuzentzailea

	1		
	Doit.	Est.	$F_1$
Hizkuntz.2: goi-muga	76,1	76,4	76,3
Ikask. autom. KPI komagabearekin + CG erregelak	77,6	52,7	62,8
<b>Ebaluazio kualitatiboa</b>	<b>83,01</b>	58,46	<b>68,61</b>

Taula: Koma zuzentzailearen emaitzak test-corpusean

# Ondorioak

- Ikasketa automatikoko teknikak (FR-Perceptron), baliagarriak
- Kateak identifikatzea, perpausak baino errazago
- Euskarako kateak eta perpausak, ingelesekoak baino zailago
- Komaren erabilera antzekoa da euskaraz eta ingelesez
- Kate- eta perpaus-identifikatzaileen informazioa esanguratsua koma-zuzentzailearentzat
- Corpusean eta hizkuntza-ezagutzan oinarritutako teknikak uztartzeak emaitzak hobetzen ditu
- Koma-zuzentzailearen ebaluaziorako ez da nahikoa aukera bakarra ontzat ematea

## Etorkizuneko lanak

- Kate- eta perpaus-identifikatzaileak hobetzea
  - Ikasketa-corpora handituz
  - Euskarako analizatzaile eta desanbiguatzailea hobetuz

## Etorkizuneko lanak

- Kate- eta perpaus-identifikatzaileak hobetzea
  - Ikasketa-corpora handituz
  - Euskarako analizatzaile eta desanbiguatzailea hobetuz
- Koma-zuzentzailea hobetzea
  - Uneko tokenaren aurreko komak kontuan hartuz
  - Emaidza-klaseen desorekaren arazoa konponduz
  - Koma okerrekin lortutako informazio linguistikoarekin gertatzen dena aztertuz



## Etorkizuneko lanak

- Kate- eta perpaus-identifikatzaileak hobetzea
  - Ikasketa-corpora handituz
  - Euskarako analizatzaile eta desanbiguatzailea hobetuz
- Koma-zuzentzailea hobetzea
  - Uneko tokenaren aurreko komak kontuan hartuz
  - Emaidza-klaseen desorekaren arazoa konponduz
  - Koma okerrekin lortutako informazio linguistikoarekin gertatzen dena aztertuz
- FR-Perceptron algoritmoa hobetzea
- Bestelako erroreen detekzioa lantzea ikasketa automatikoko tekniken bidez
- Analisi-katean eta XUXENg-n integratzea

## Tesiari hertsiki lotutako argitalpenak (I)

- Alegria I., Arrieta B., Carreras X., Díaz de Ilarraza A., Uria L. **Chunk and Clause Identification for Basque by Filtering and Ranking with Perceptrons.** *Revista del procesamiento del lenguaje natural*, nº 41, pags: 5-12; 2008.
- Aldabe I., Arrieta B., Díaz de Ilarraza A., Maritxalar M., Niebla I., Oronoz M., Uria L. **Basque error corpora: a framework to classify and store it.** *In the Proceedings of the 4th Corpus Linguistic Conference.* Birmingham. UK. 2007.
- Alegria I., Arrieta B., Díaz de Ilarraza A., Izagirre E., Maritxalar M. **Using Machine Learning Techniques to Build a Comma Checker for Basque.** *Proceedings of Coling-ACL.* Sydney. Australia. 2006.
- Aduriz I., Arrieta B., Arriola J.M., Díaz de Ilarraza A., Izagirre E., Ondarra A. **Muga Gramatikaren Optimizazioa.** *EHU/LSI/TR 26-2005.* Donostia. Euskal Herria. 2005
- Aldabe I., Arrieta B., Díaz de Ilarraza A., Maritxalar M., Oronoz M., Uria L. **Propuesta de una clasificación general y dinámica para la definición de errores.** *Revista de Psicodidáctica, EHU. Vol 10, Nº 2, p. 47-60.* 2005.

## Tesiari hertsiki lotutako argitalpenak (II)

- Ansa O., Arregi X., Arrieta B., Díaz de Ilarraza A., Ezeiza N., Fernandez I., Garmendia A., Gojenola K., Laskurain B., Martínez E., Oronoz M., Otegi A., Sarasola K., Uria L. **Integrating NLP Tools for Basque in Text Editors.** *Workshop on International Proofing Tools and Language Technologies.* University of Patras. Greece. 2004.
- Aldezabal I., Aranzabe M., Arrieta B., Maritxalar M., Oronoz M. **Toward a punctuation checker for Basque.** *ATALA workshop. Le role de la typographie et de la ponctuation dans le traitement automatique des langues.* Paris. France. 2003.
- Arrieta B., Díaz de Ilarraza A., Gojenola K., Maritxalar M., Oronoz M. **A database system for storing second language learner corpora.** *Learner corpora workshop. Corpus linguistics 2003. Volume 16, Part 1. p.: 33-41;* Lancaster, UK. 2003.
- Aduriz I., Aldezabal I., Aranzabe M., Arrieta B., Arriola J., Atutxa A., Díaz de Ilarraza A., Gojenola K., Oronoz M., Sarasola K., Urizar R. **The design of a digital resource to store the knowledge of linguistic errors.** *DRH2002 (Digital Resources for the Humanities), pp 76-78.* Edinburgh. Scotland. 2002.

*“Euskaran bezala bizitzan, gauza ustez txiki baina guztiz beharrezkoak baztertzeko joera dugu eta nahiago ditugu proiektu pinpirinak. Baina elementalena —eta neketsuena— koma baten garrantzia eta gisakoak azpimarratzea da.”*

A. Lertxundi (Berria; 2010/07/15)

Eskerrik asko!

Azaleko sintaxiaren tratamendua  
ikasketa automatikoko tekniken bidez:  
euskarako kateen eta perpausen identifikazioa  
eta bere erabilera koma-zuzentzaile batean

**Doktoregaia:** Bertol Arrieta Kortajarena

**Zuzendariak:** Iñaki Alegria Loinaz

Arantza Diaz de Ilarraza Sanchez

Lengoaia eta Sistema Informatikoak/Lenguajes y Sistemas Informáticos  
Euskal Herriko Unibertsitatea/Universidad del País Vasco

2010eko uztailaren 27a