

Tuning As Ranking

Mark Hopkins
Jonathan May
SDL Language Weaver

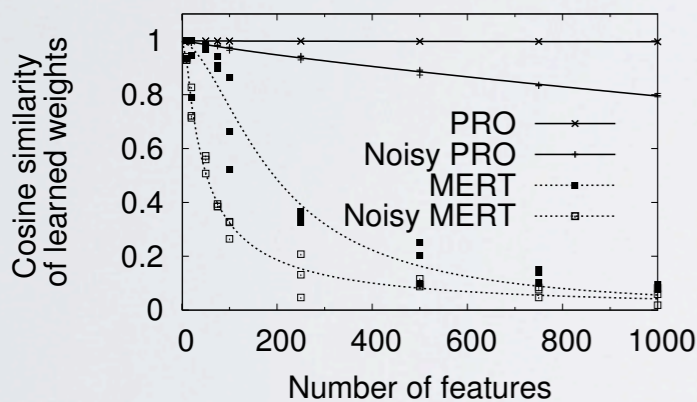
EMNLP
July 29, 2011

What we did

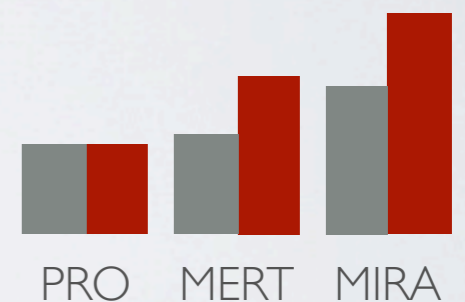
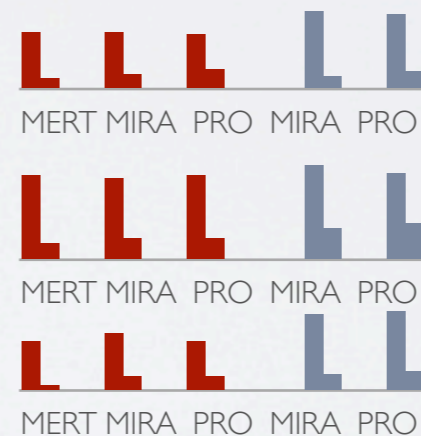
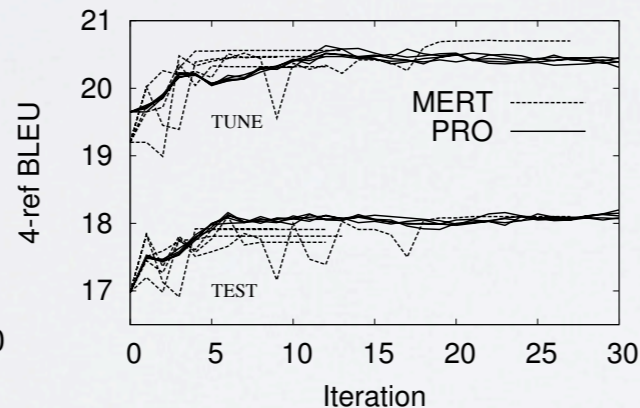
We replaced **MERT's linear optimization** with a **linear binary classifier**, and fed it **pairs** of translations, effecting a **ranking**

What we found

Synthetic weight learning of MERT and PRO



Urdu-English PBMT tuning stability



Scalable to many features

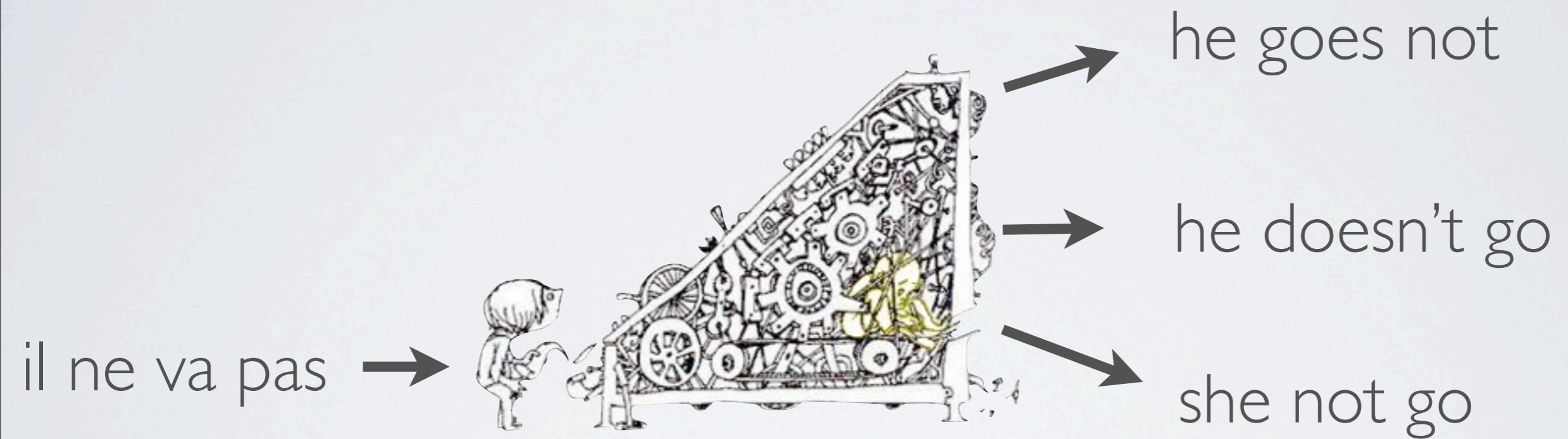
Consistent results

Parity with leading techniques

Very **fast**

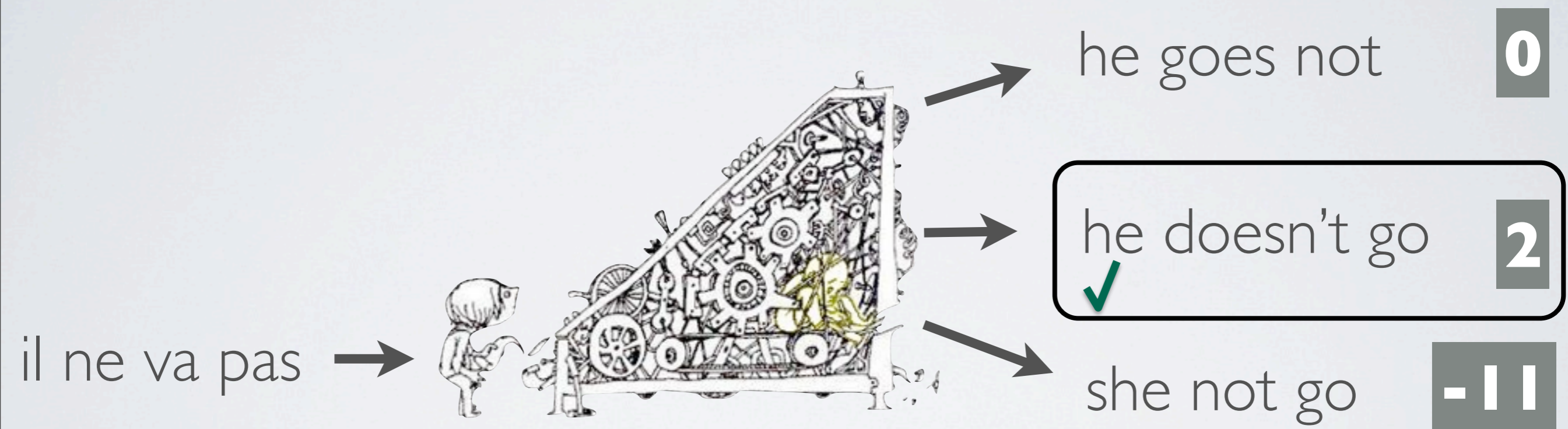
Any Questions?

Which is best?



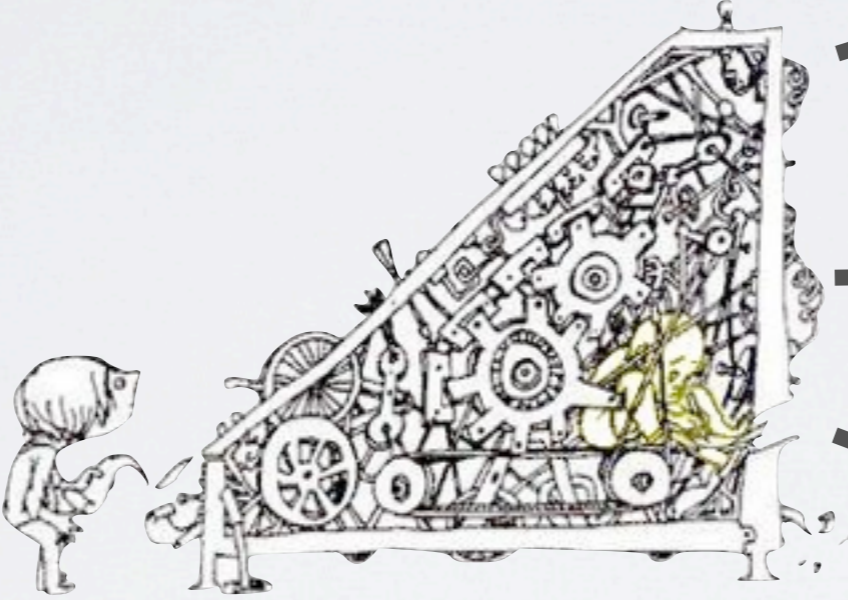
(Image credit: Silverstein, 1981)

Which is best?



A good scoring function can tell us

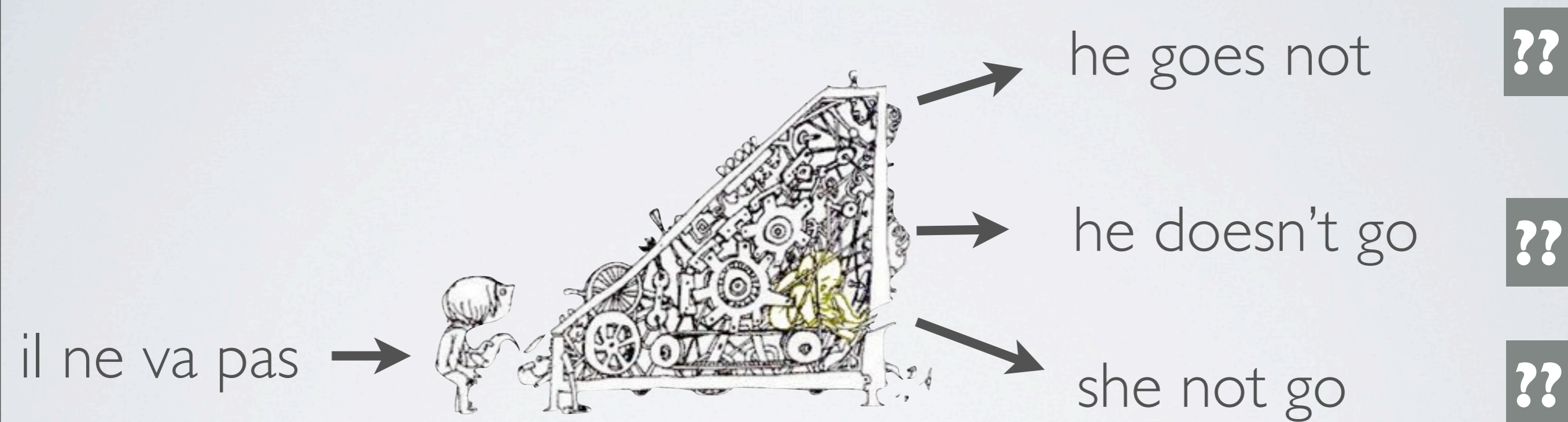
Which is best?

il ne va pas → 

he goes not	-
he doesn't go	0
she not go	2

We should avoid bad functions

Which is best?



How do we ensure “proper” scores?

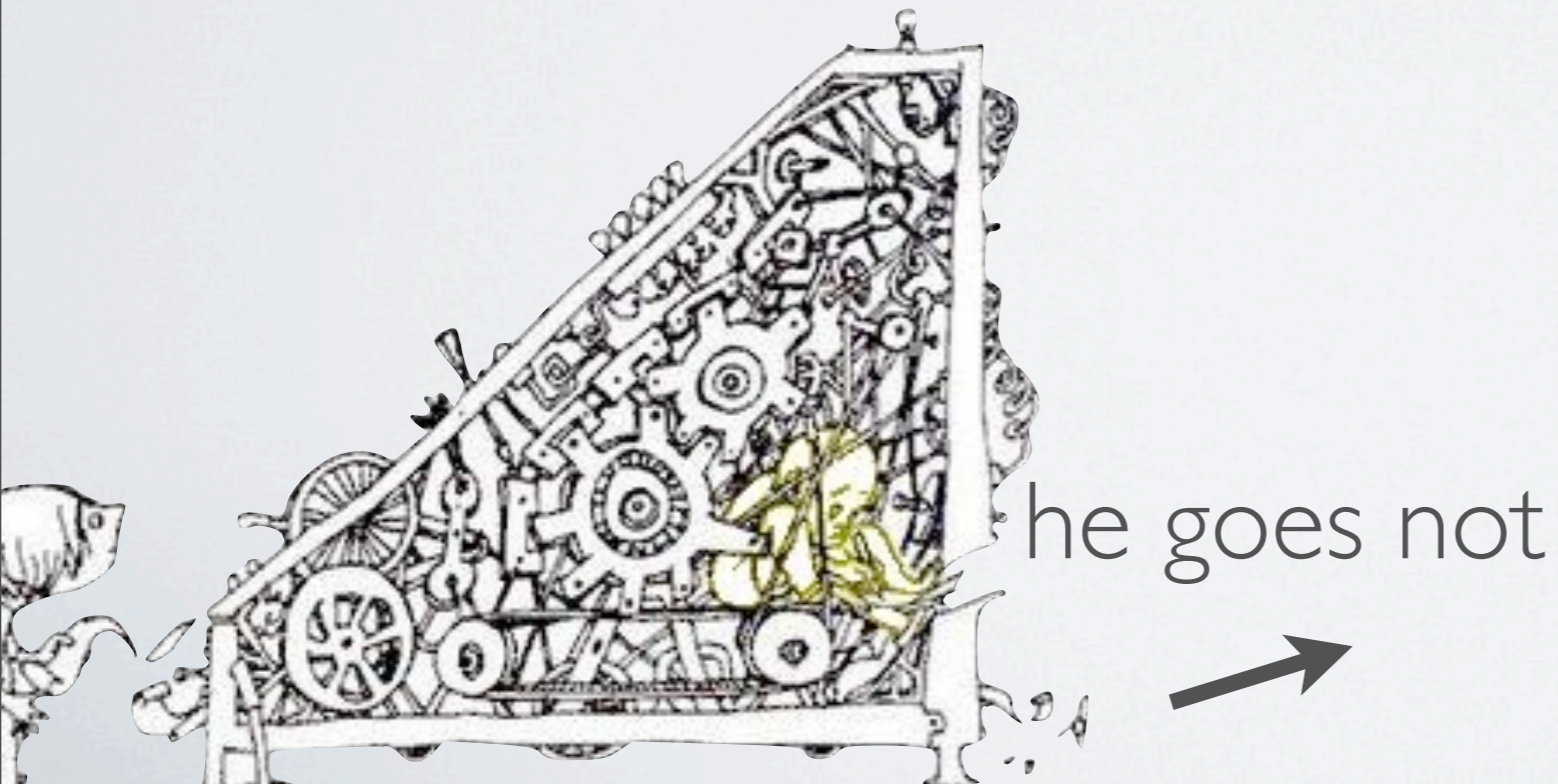
Properties of the translation



Properties of the translation

literal meaning?

2

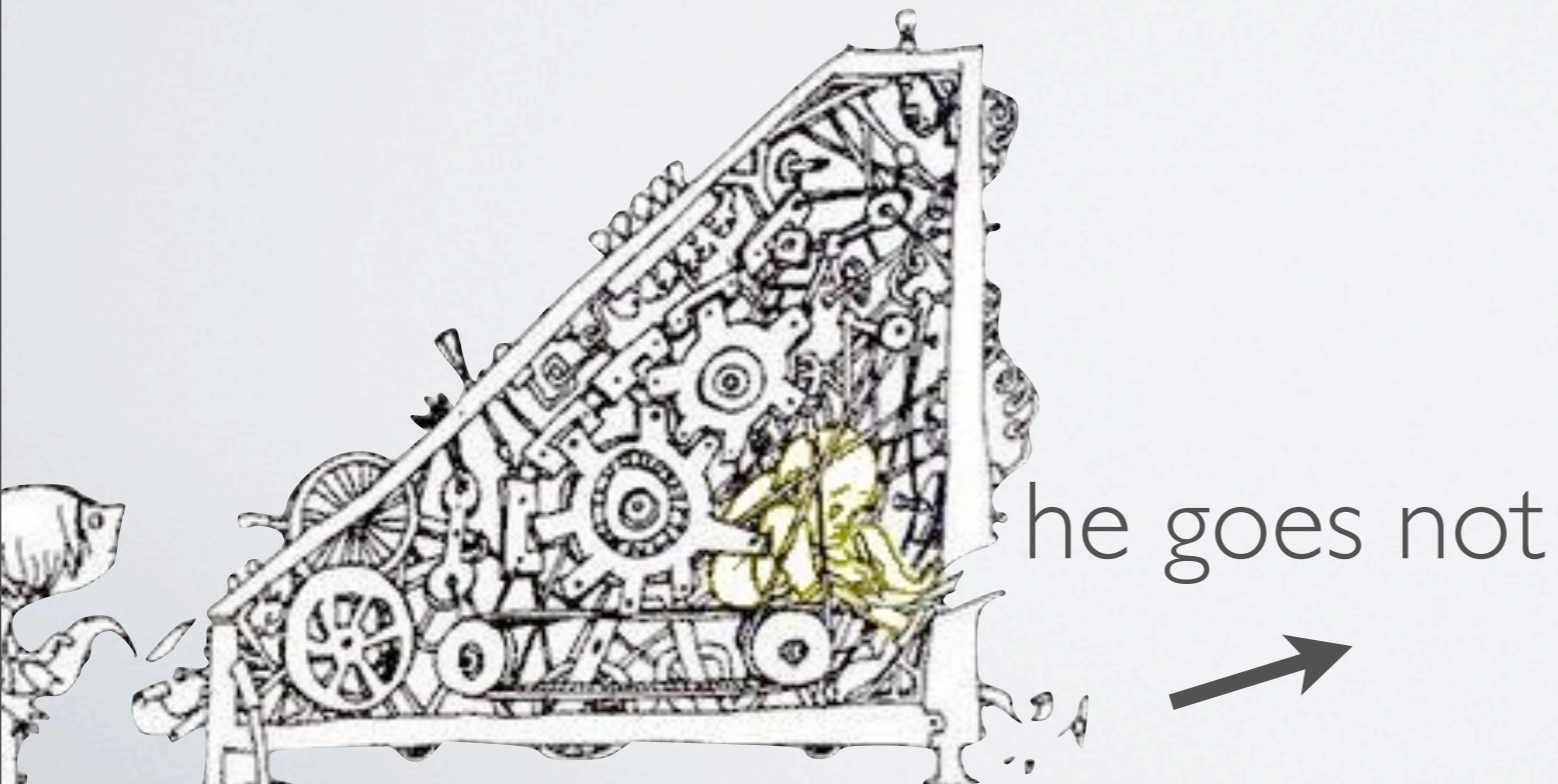


Properties of the translation

literal meaning?
fluency?

2

4



Properties of the translation

literal meaning?

2

fluency?

4

word count?

3

count of "he"?

1

count of "coffee"?

0

...

alliterative?

0

$2 \leq \text{count}(\text{"o"}) \leq 3$?

1

how do you feel? **32.8**



Properties of the translation

Features!

literal meaning?

2

fluency?

4

word count?

3

count of "he"?

1

count of "coffee"?

0

...

alliterative?

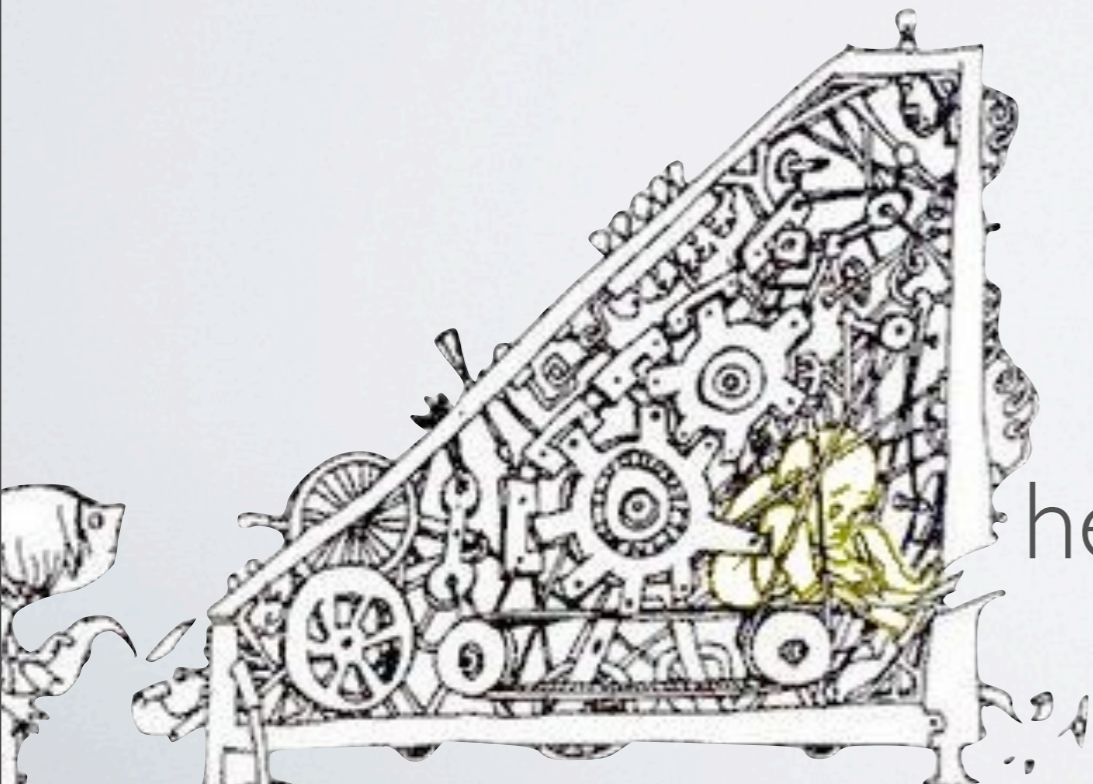
0

$2 \leq \text{count}(\text{"o"}) \leq 3$?

1

how do you feel? **32.8**

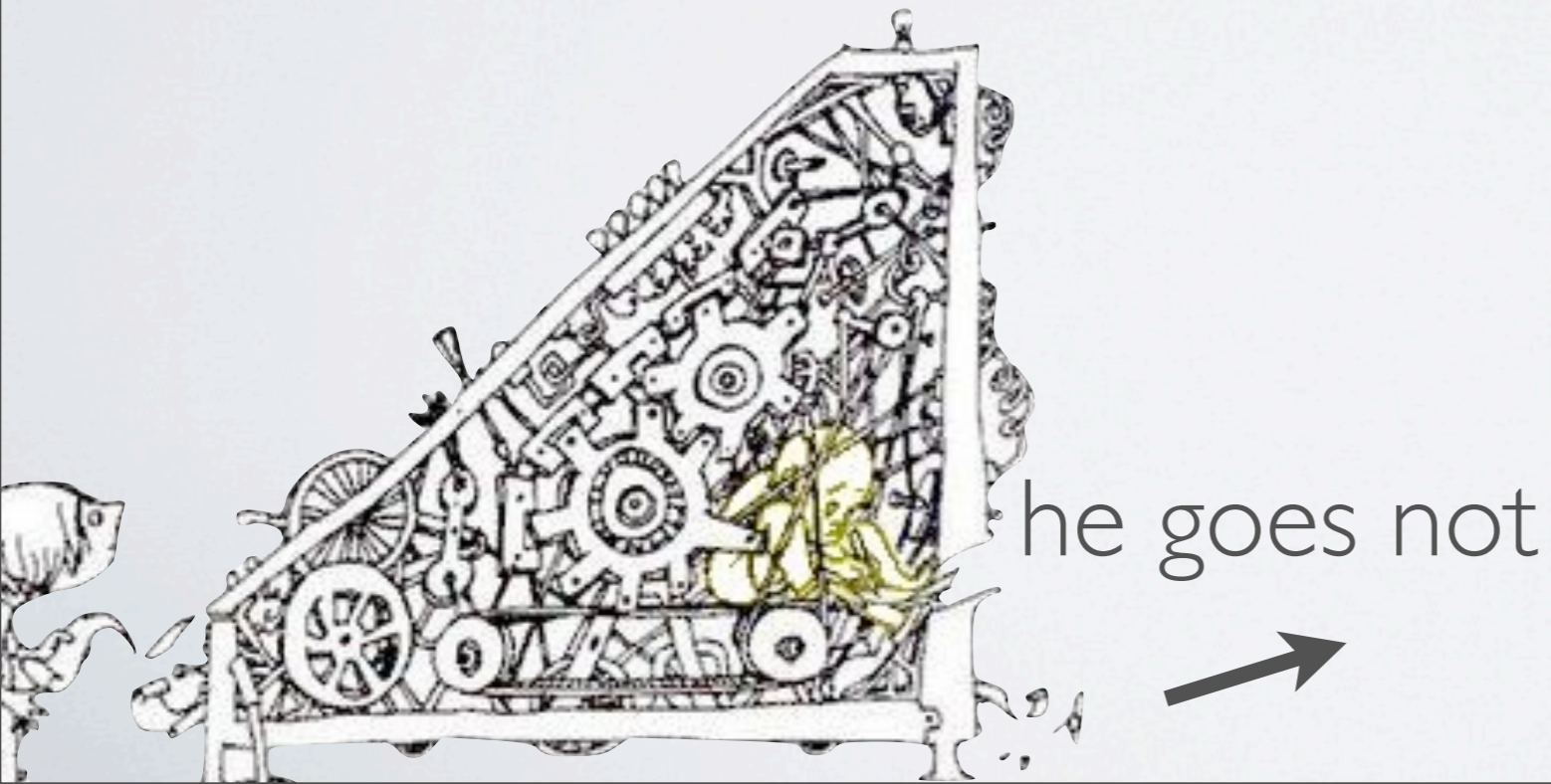
he goes not



Properties of the translation

Features!

f1: 2
f2: 4



Properties of the translation

Features!

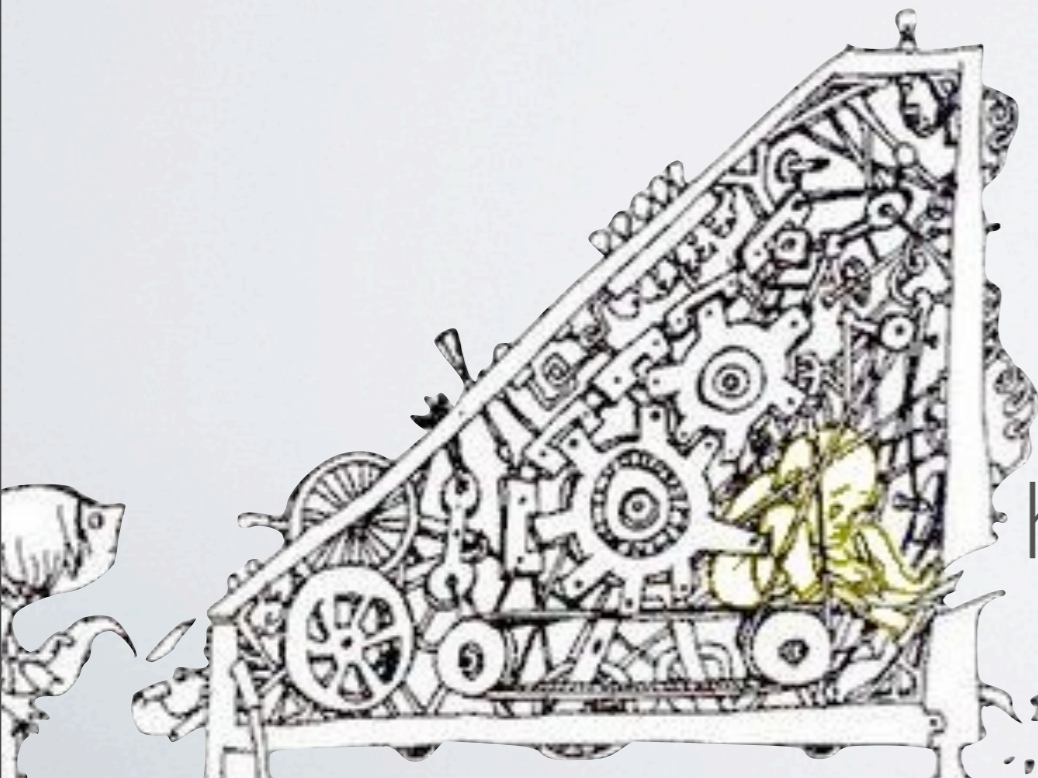
f1:
f2:

Form a weighted sum

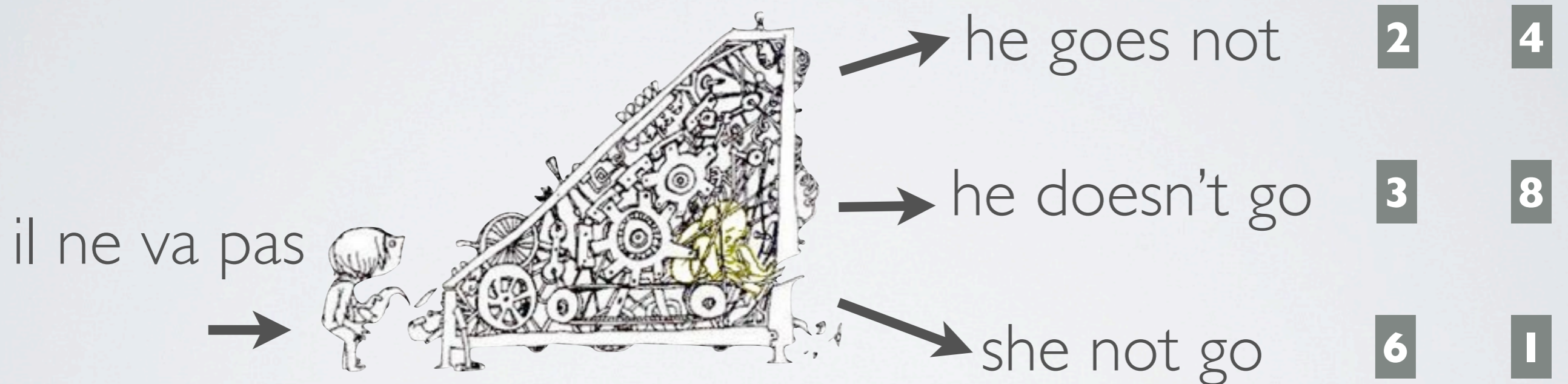
$$\mathbf{-2} \times \mathbf{2} + \mathbf{3} \times \mathbf{4} = \mathbf{8}$$

↑ ↑
Weights!

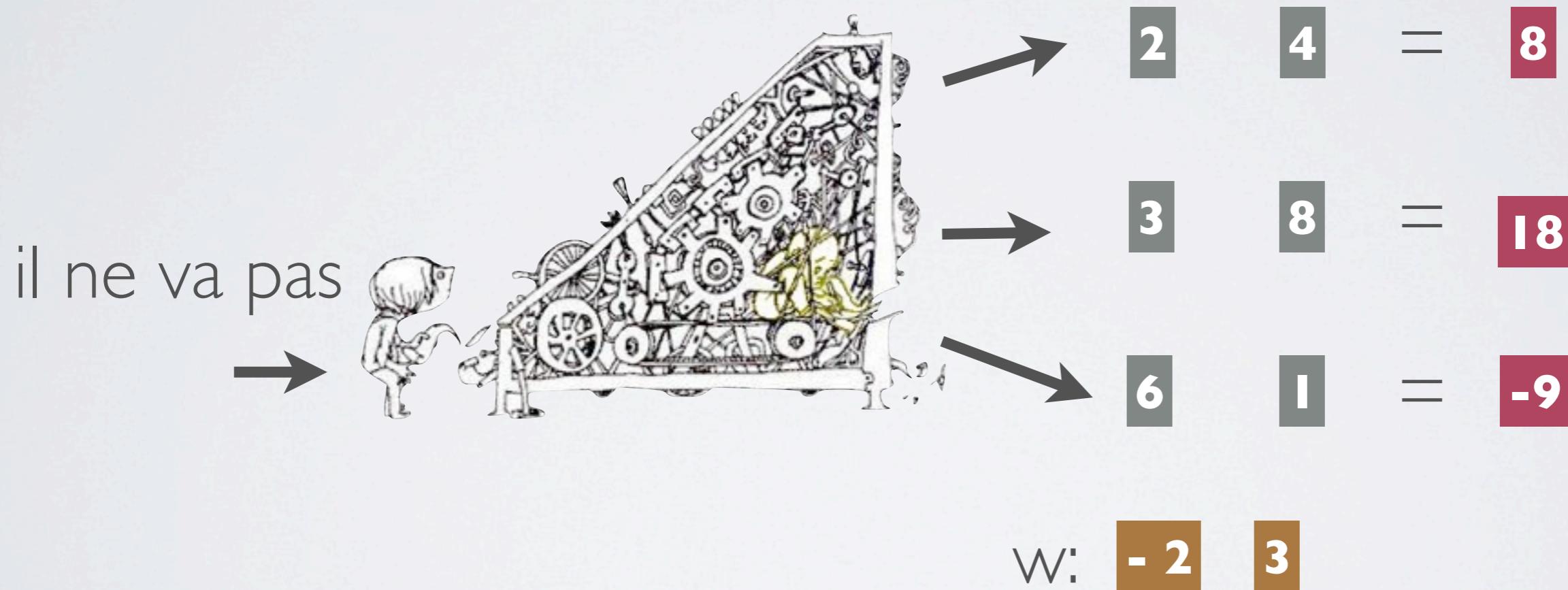
he goes not



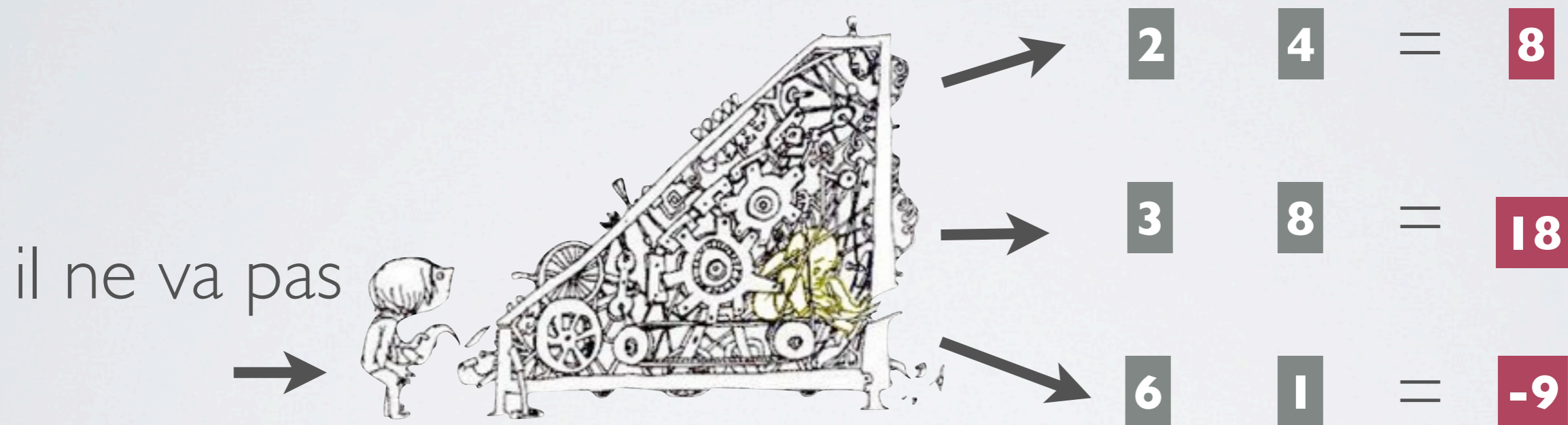
Translations are feature vectors



Weight vector determines the score



Weight vector determines the score

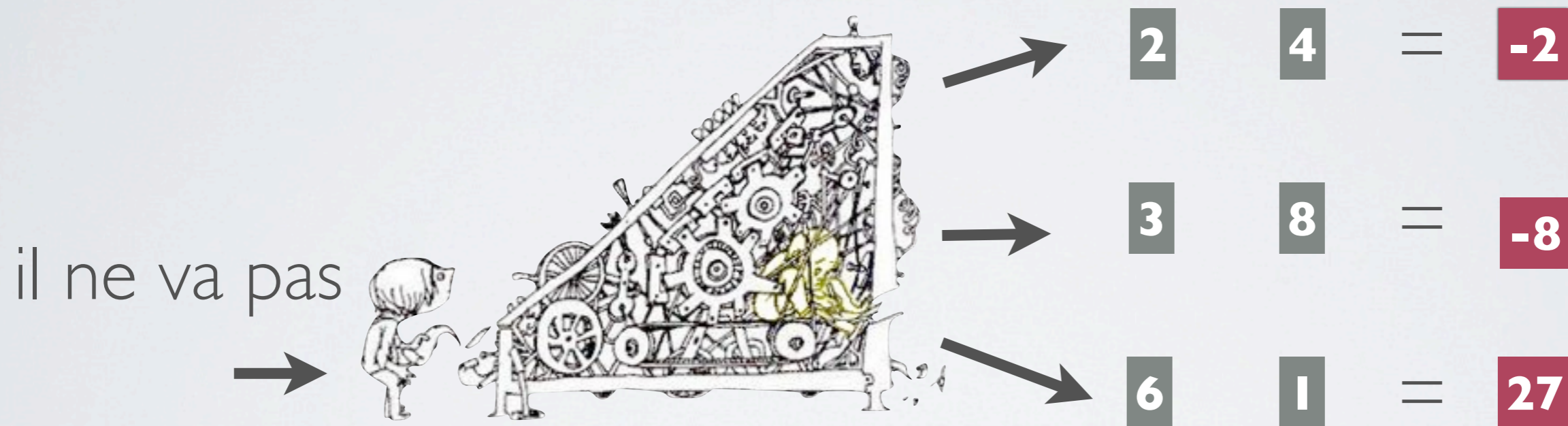


model score = features • weights

$h = f \cdot w$

W: -2 3

Weight vector determines the score

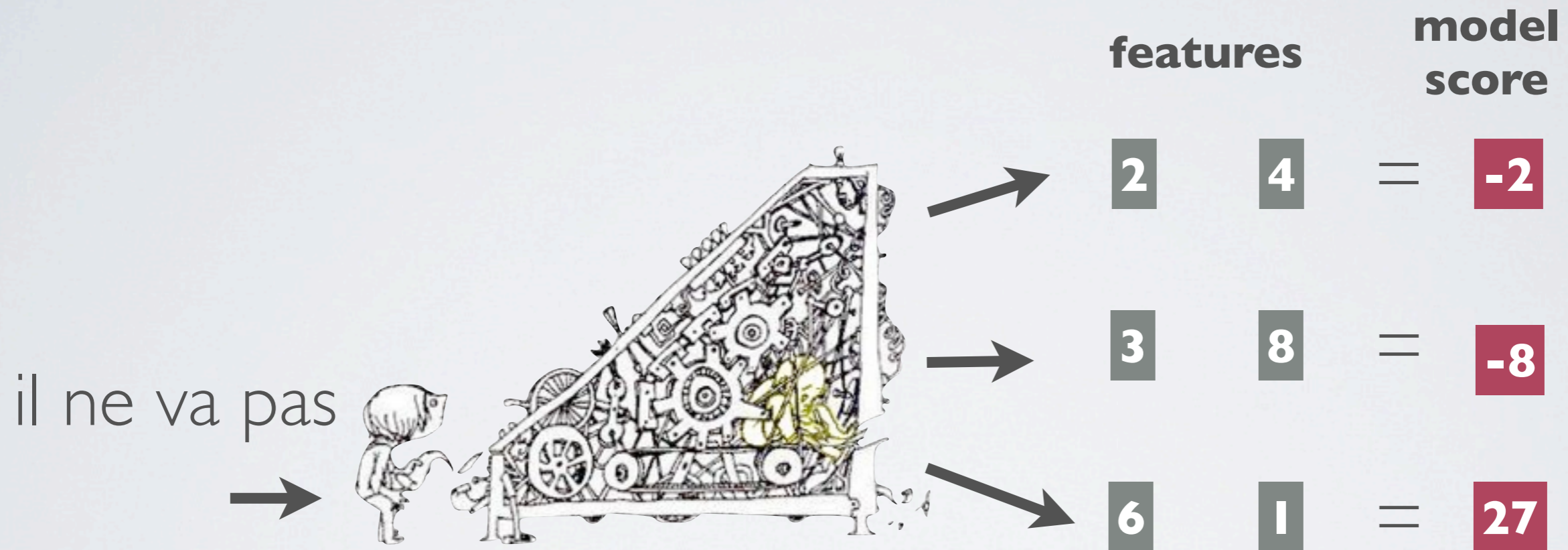


model score = features • weights

$h = f \cdot w$

W: 5 -3

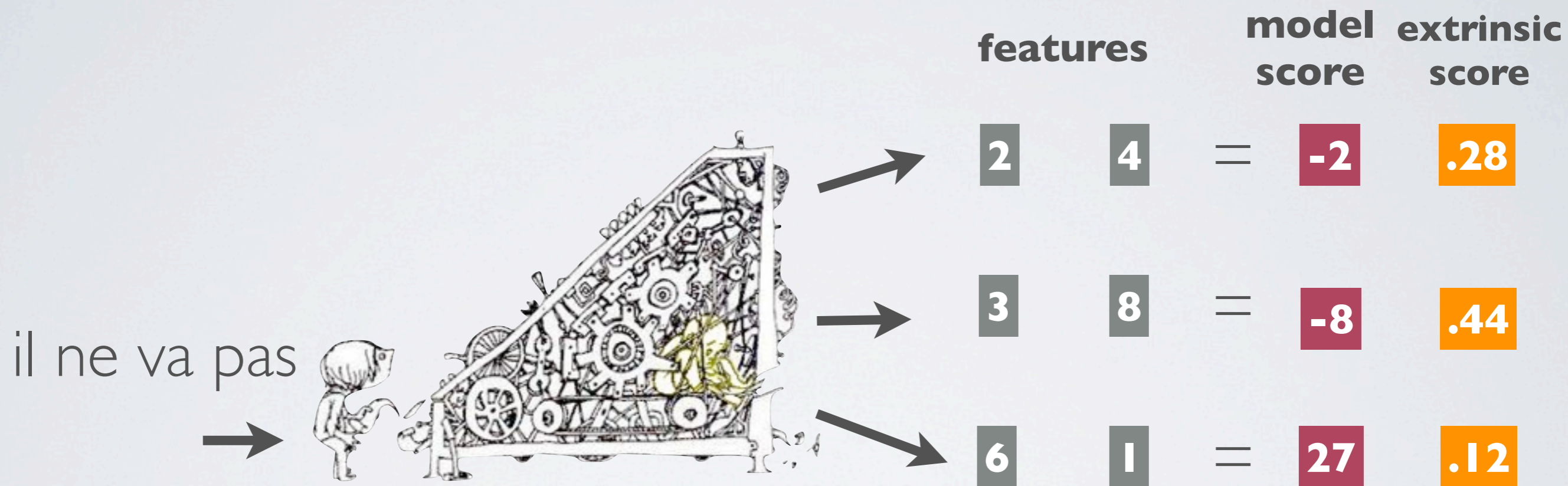
Weight vector determines the score



Tuning is all about choosing this vector

w : 5 -3

Weight vector determines the score

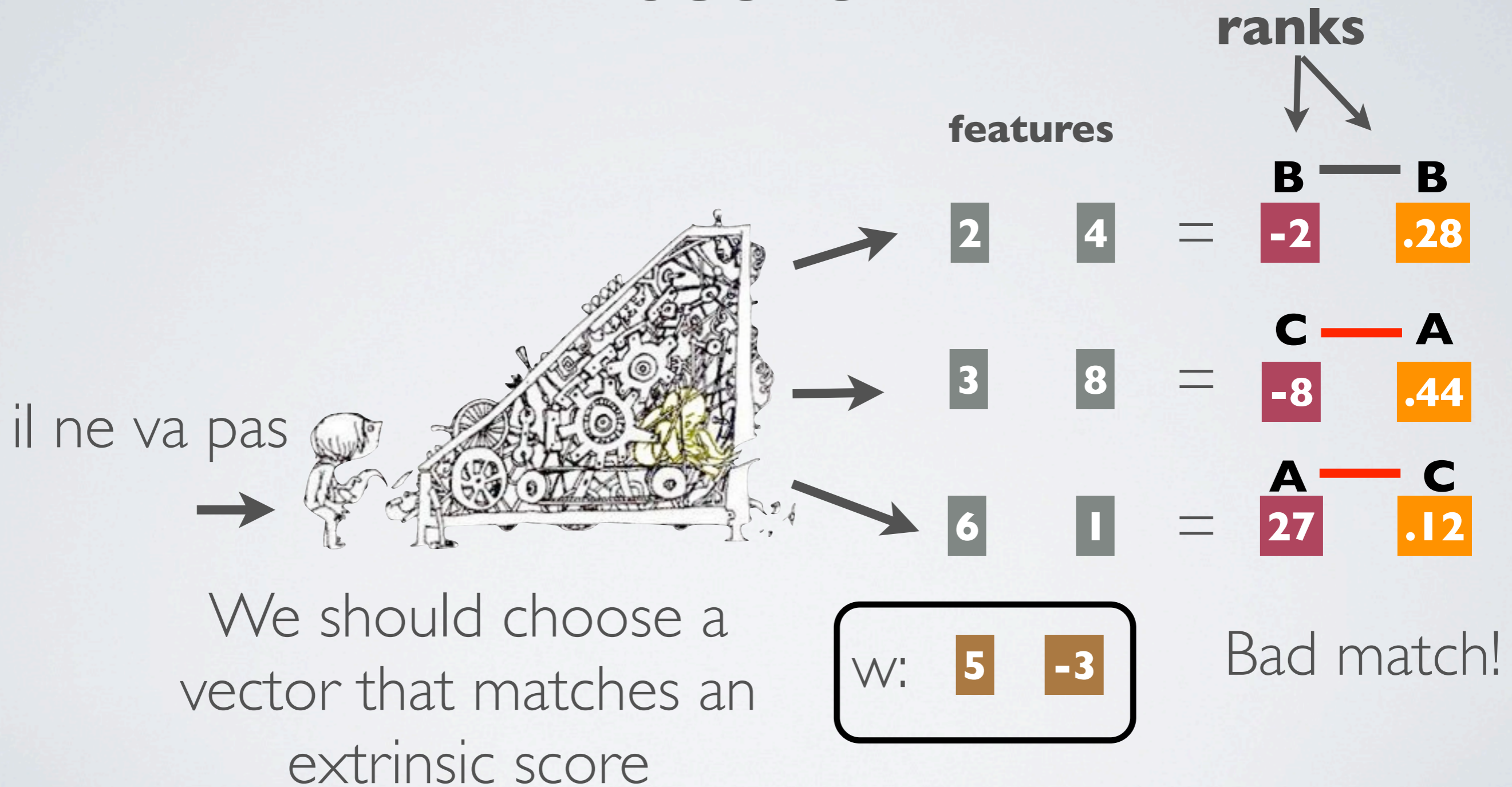


il ne va pas

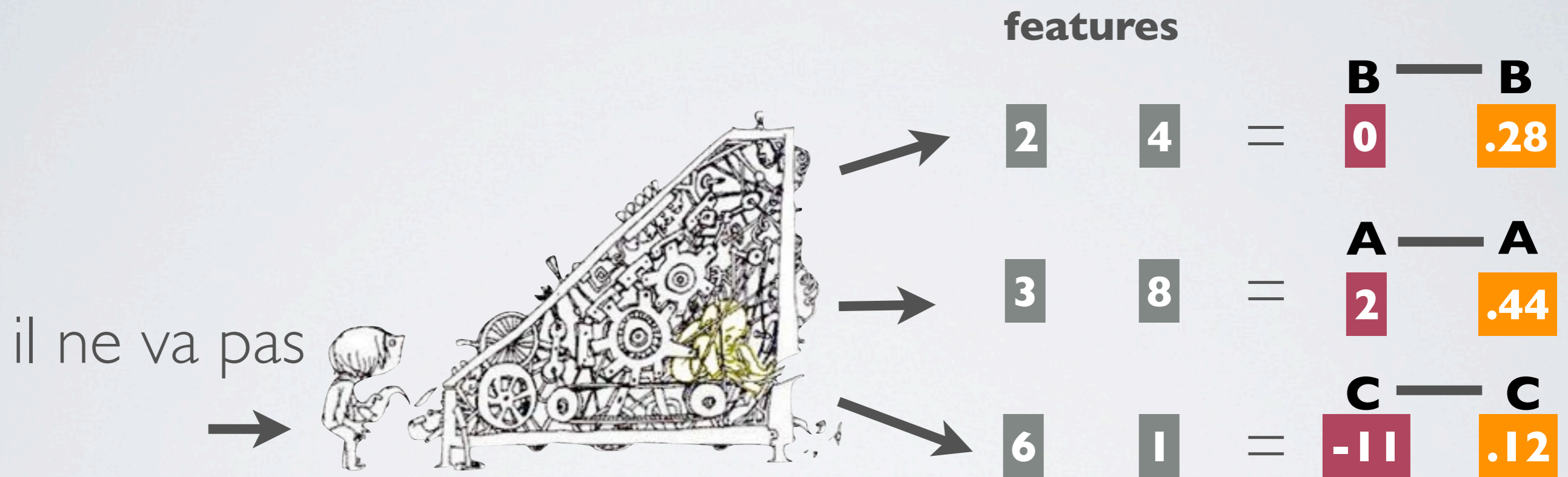
We should choose a vector that matches an extrinsic score

(Lin & Och, '04)

Weight vector determines the score



Weight vector determines the score



We should choose a vector that matches an extrinsic score

w: **-2** **1**

Good match!

The tuning framework that
everybody uses

MERT framework

(Och, 2003)

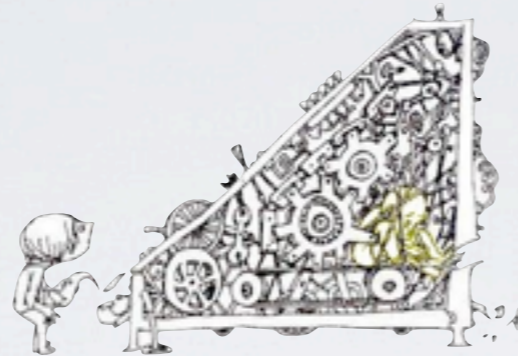
The tuning framework that
*(almost)** everybody uses

MERT framework

** Not David Chiang*

(Och, 2003)

The tuning framework that *(almost)** everybody uses



**Candidate
Generation**

MERT framework

Generate n-best
per input sentence

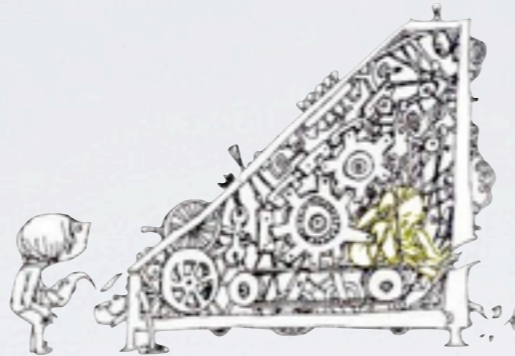


**Add to
Candidate Pool**

** Not David Chiang*

(Och, 2003)

The tuning framework that *(almost)** everybody uses



**Candidate
Generation**

Generate n-best
per input sentence

MERT framework



**Add to
Candidate Pool**



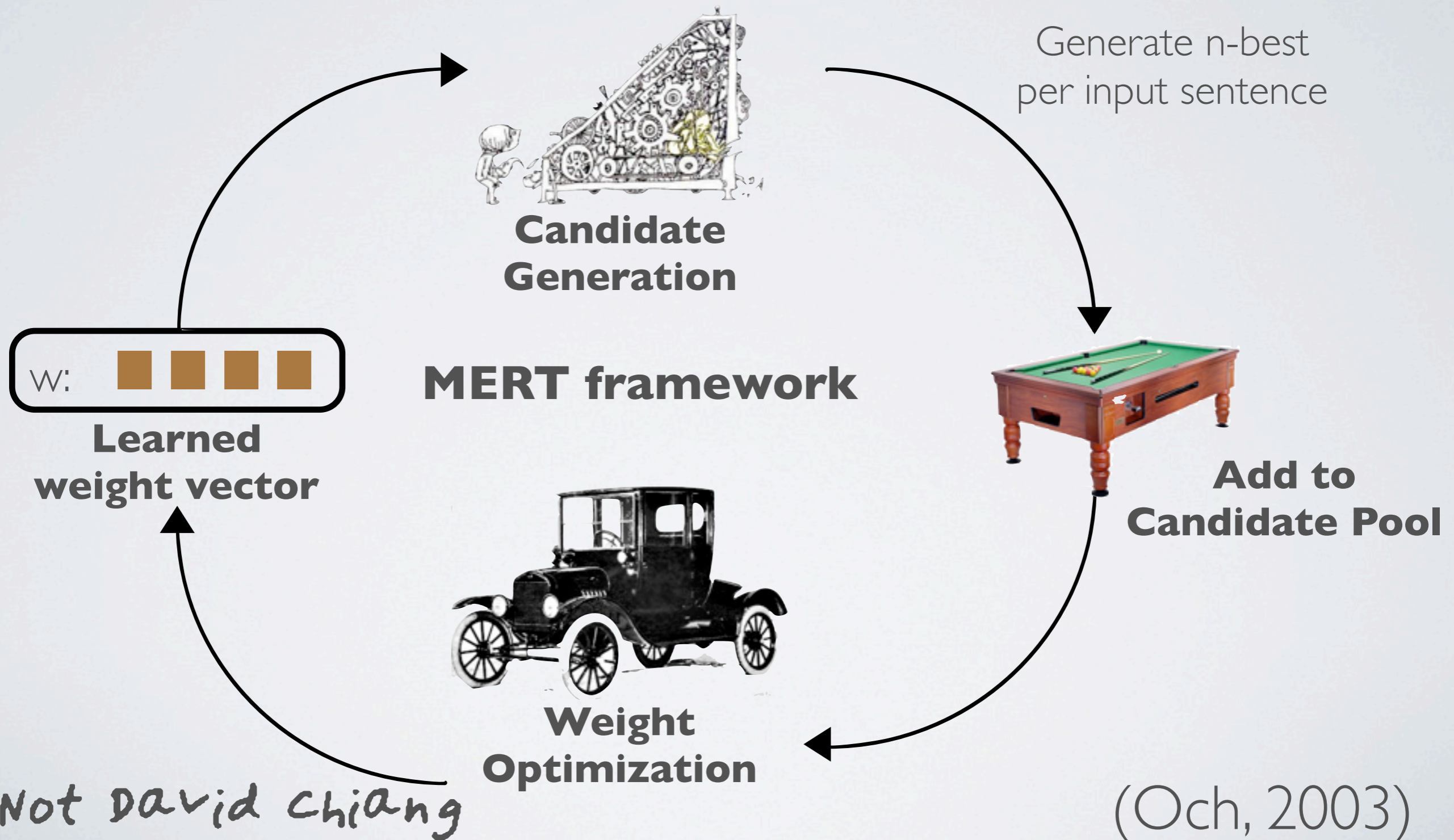
**Weight
Optimization**



** Not David Chiang*

(Och, 2003)

The tuning framework that *(almost)** everybody uses



How MERT works



S1

S2

feats

extrins

2	4
3	8
6	1

-3	-3
1	5
-5	-3

.28
.44
.12
.15
.18
.32

The MERT algorithm works by varying one weight at a time to find a value for that weight that aligns the **best model score** with the **best extrinsic score**



How MERT works



S1

feats model extrins

2	4	0	.28
3	8	0	.44
6	1	0	.12

S2

-3	-3	0	.15
1	5	0	.18
-5	-3	0	.32

W: 0 0

total extrinsic **.43**

The MERT algorithm works by varying one weight at a time to find a value for that weight that aligns the **best model score** with the **best extrinsic score**

How MERT works



S1

S2

feats model extrins

2	4	0	.28
3	8	0	.44
6	1	0	.12
-3	-3	0	.15
1	5	0	.18
-5	-3	0	.32

W: 0 0

total **.43**
extrinsic

(MERT can optimize the non-decomposable BLEU; swap these for n-gram component values and determine **total** with the BLEU equation)

How MERT works



S1

feats model extrins

2	4	0	.28
3	8	0	.44
6	1	0	.12

S2

-3	-3	0	.15
1	5	0	.18
-5	-3	0	.32

W: 0 0

↑
vary

↑
hold

**total
extrinsic** **.43**

The MERT algorithm works by varying one weight at a time to find a value for that weight that aligns the **best model score** with the **best extrinsic score**

How MERT works



	feats	model	extrins
S1	2	4	-2
	3	8	-3
	6	1	-6
S2	-3	-3	3
	1	5	-1
	-5	-3	5
			total extrinsic
			.60

W:	-1	0
	↑ vary	↑ hold

The MERT algorithm works by varying one weight at a time to find a value for that weight that aligns the **best model score** with the **best extrinsic score**

How MERT works



S1

feats model extrins

2	4	-2	.28
3	8	-3	.44
6	1	-6	.12

S2

-3	-3	3	.15
1	5	-1	.18
-5	-3	5	.32

W: -1 0

↑ hold ↑ vary

total extrinsic **.60**

The MERT algorithm works by varying one weight at a time to find a value for that weight that aligns the **best model score** with the **best extrinsic score**

How MERT works



S1

S2

feats model extrins

2	4	2	.28
3	8	5	.44
6	1	-5	.12
-3	-3	0	.15
1	5	4	.18
-5	-3	2	.32

W: -1 1

↑ ↑
hold **vary**

total extrinsic .62

The MERT algorithm works by varying one weight at a time to find a value for that weight that aligns the **best model score** with the **best extrinsic score**

How MERT works



	feats	model	extrins	
S1	2	4	0	.28
	3	8	2	.44
	6	1	-11	.12
S2	-3	-3	3	.15
	1	5	1	.18
	-5	-3	7	.32

The MERT algorithm works by varying one weight at a time to find a value for that weight that aligns the **best model score** with the **best extrinsic score**



total extrinsic **.76**

How MERT works



S1

S2

feats

model

extrins

	2	4	0	2								??	.28
	3	8	1	3								??	.44
	6	1	1	-5								??	.12
	-3	-3	0	7								??	.15
	1	5	0	17								??	.18
	-5	-3	2	6								??	.32

W: **-2** **1** **3** **16**

This works well for small feature sets, but as the feature space grows, it is hard to find a good position

Synthetic Experiment

random

random

random

“features”

“Candidate pool” of randomly drawn “feature” vectors

Synthetic Experiment



“Candidate pool” of randomly drawn “feature” vectors

How to determine “extrinsic score”?

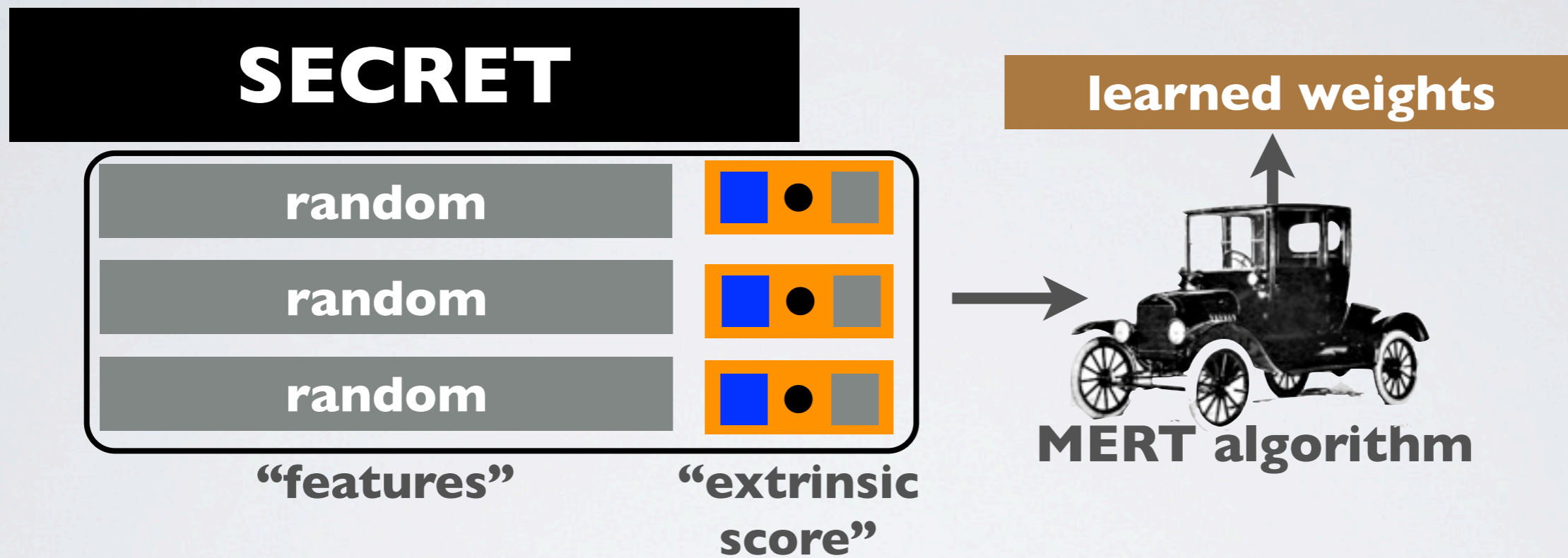
Synthetic Experiment



“Candidate pool” of randomly drawn “feature” vectors

Secret “goal weights” used to calculate extrinsic score

Synthetic Experiment

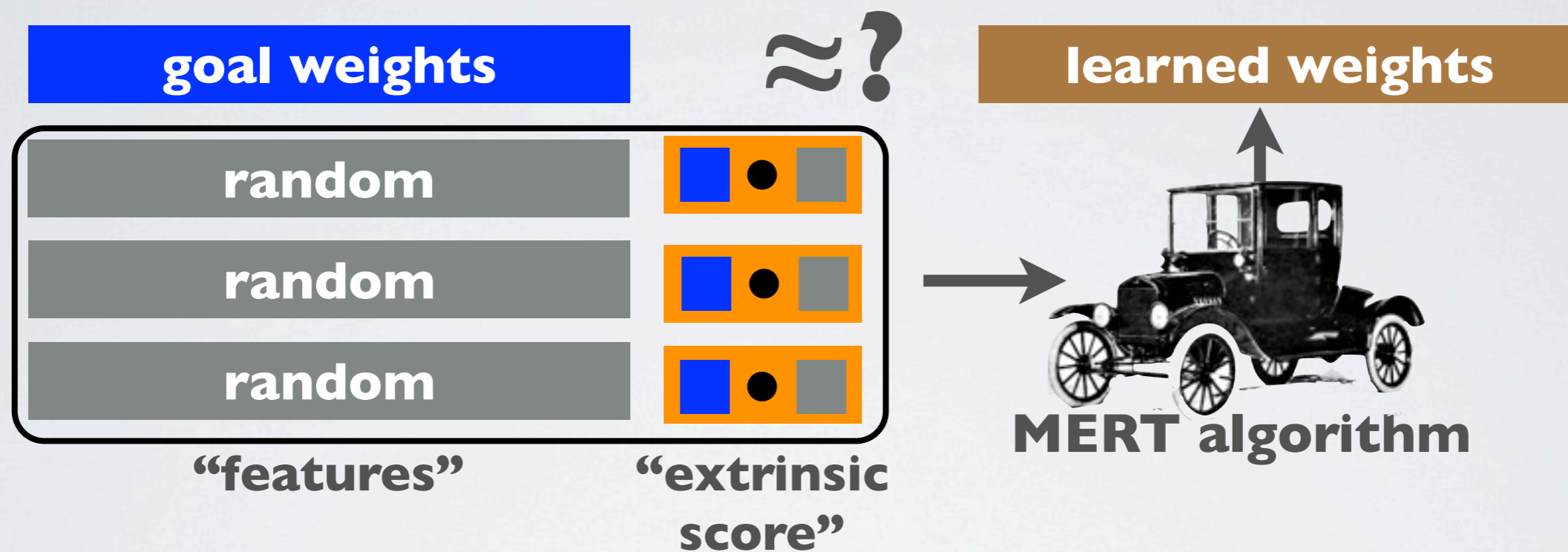


Now use MERT to try and learn the goal weights back

This is linear equation solving

It's much easier than MT tuning

Synthetic Experiment



Now use MERT to try and learn the goal weights back

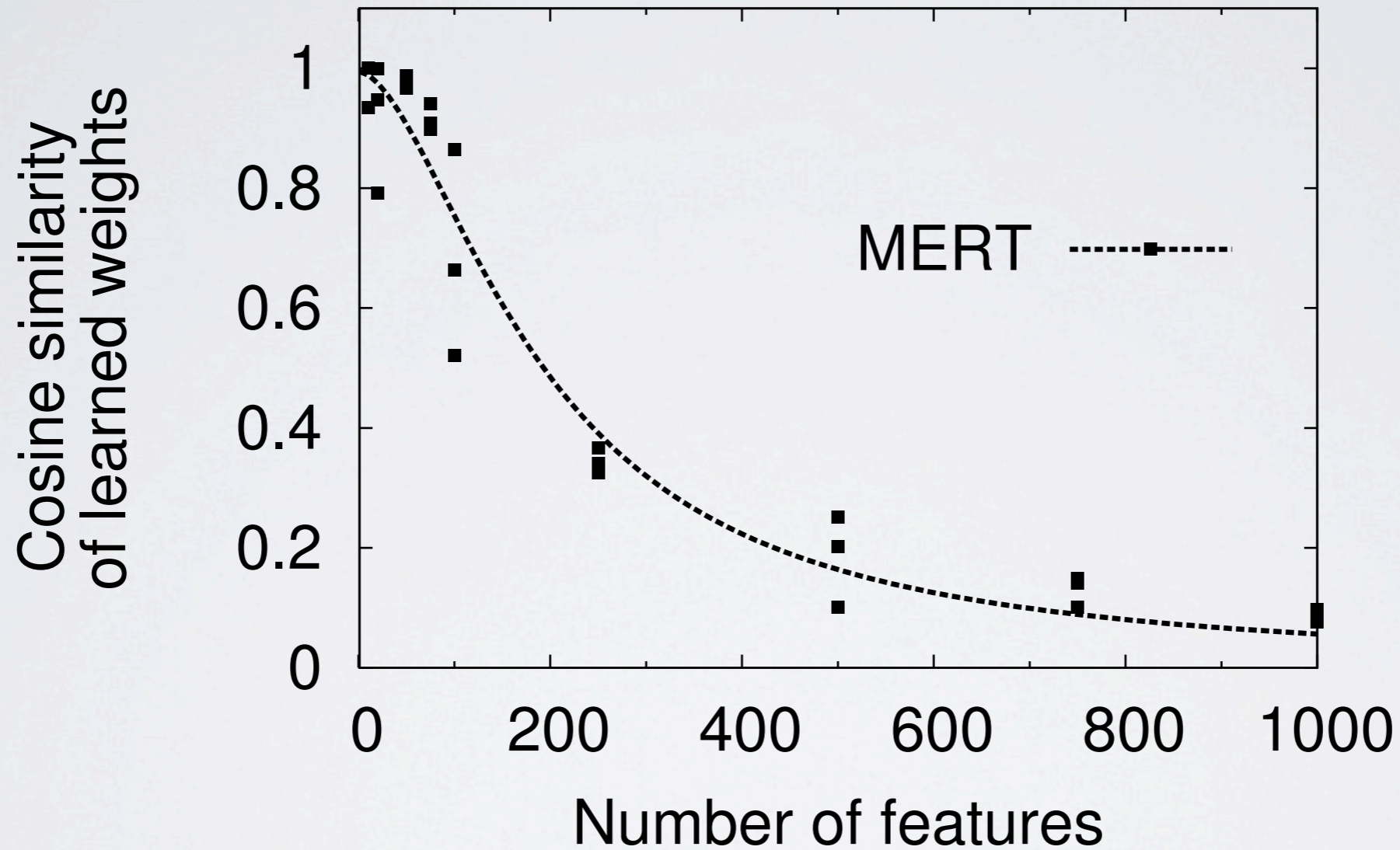
This is linear equation solving

It's much easier than MT tuning



MERT *doesn't scale*

Synthetic weight learning of MERT



The synthetic experiment in ideal conditions validates what has long been accepted as truth

MERT only cares about the top-scoring translation

feats model extrins



S1

2	4	0	B
3	8	2	A
6	1	-11	C

S2

-3	-3	3	C
1	5	1	B
-5	-3	7	A



MERT only cares about the top-scoring translation

feats model extrins

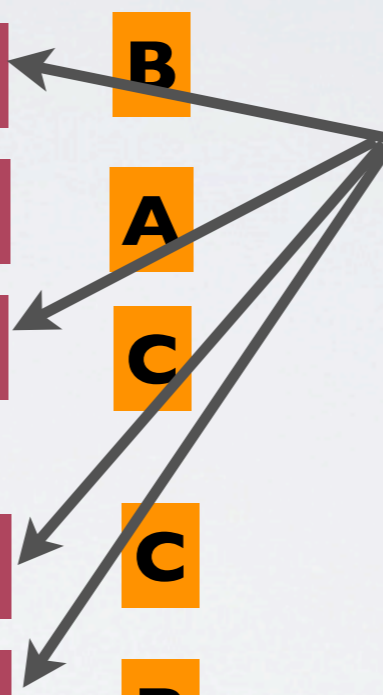


S1

S2

2	4	??	B
3	8	A	A
6	1	??	C
-3	-3	??	C
1	5	??	B
-5	-3	A	A

MERT doesn't care about these



It doesn't care about matching the overall ranking

feats model extrins



S1

2	4	B	B
3	8	A	A
6	1	C	C

S2

-3	-3	B	C
1	5	C	B
-5	-3	A	A

mismatch!

This could lead to poor generalization

feats model extrins



2	4	D	B
3	8	A	A
6	1	C	C
7	-4	J	S
12	-2	B	G
...



not good but
liked by model

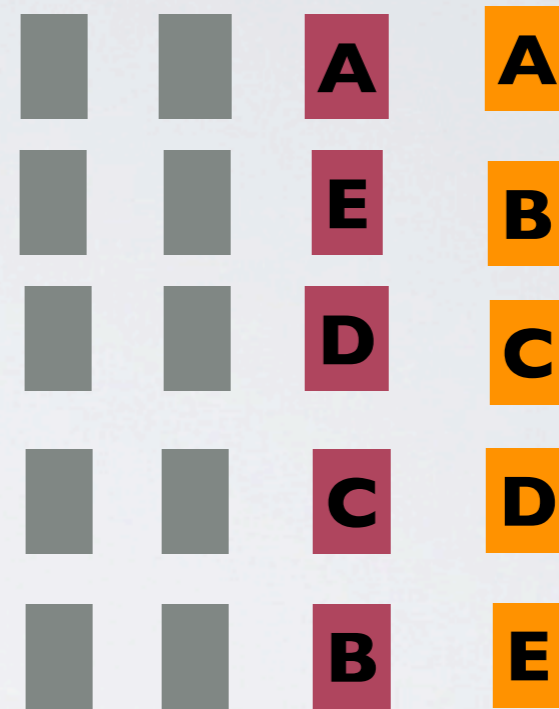
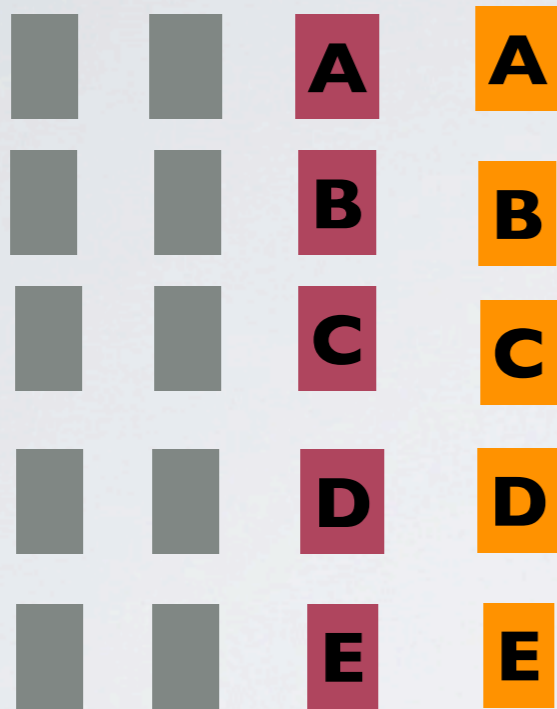
SI

-12	5	Z	D
-----	---	---	---



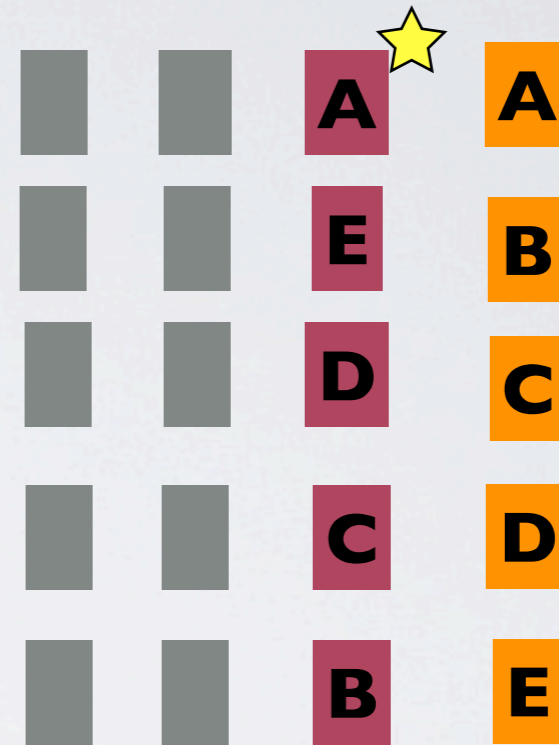
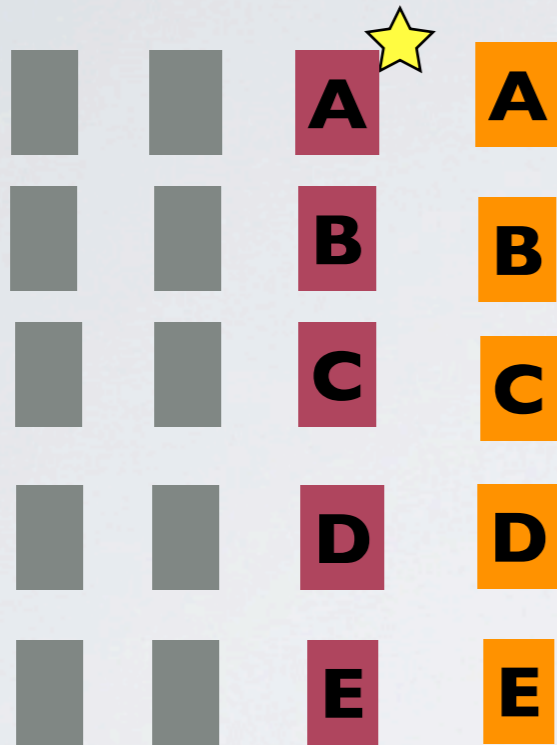
good but
disliked

We should focus on rank



Recognize that these
are different solutions!

We should focus on rank



Recognize that these
are different solutions!
(To MERT they are the
same)

We can describe rank from a
pairwise perspective

translation a

\vec{f}_a h_a g_a

translation b

\vec{f}_b h_b g_b

For any two translations **a** and **b**
of the same sentence

(Herbrich et al., '99)

We can describe rank from a
pairwise perspective

translation a

\vec{f}_a h_a g_a

translation b

\vec{f}_b h_b g_b

extrinsic

$g_a > g_b$



model

$h_a > h_b$

Model and extrinsic
score **order** should agree

We can describe rank from a **pairwise** perspective

translation a

\vec{f}_a h_a g_a

translation b

\vec{f}_b h_b g_b

extrinsic

$g_a > g_b$



model

$h_a > h_b$

$h_a - h_b > 0$

We can describe rank from a **pairwise** perspective

translation a



translation b



extrinsic

$$g_a > g_b$$



model

$$h_a > h_b$$

$$h_a - h_b > 0$$

$$\vec{w} \cdot \vec{f}_a - \vec{w} \cdot \vec{f}_b > 0$$

We can describe rank from a **pairwise** perspective

translation a



translation b



extrinsic

$$\vec{g}_a > \vec{g}_b$$



model

$$\vec{h}_a > \vec{h}_b$$

$$\vec{h}_a - \vec{h}_b > 0$$

$$\vec{w} \cdot \vec{f}_a - \vec{w} \cdot \vec{f}_b > 0$$

$$\vec{w} \cdot \vec{f}_a - \vec{f}_b > 0$$

This is a **binary classification** problem

extrinsic

$$g_a > g_b$$



model

$$\vec{w} \cdot \vec{f_a - f_b} > 0$$

This is a **binary classification** problem

extrinsic

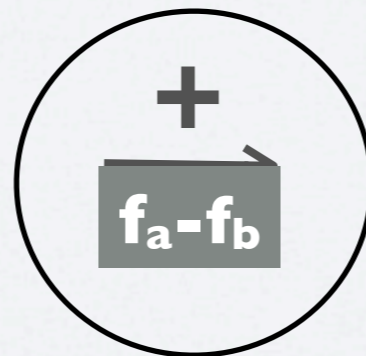
$$g_a > g_b$$

label
(+ if a is better,
- if b is better)

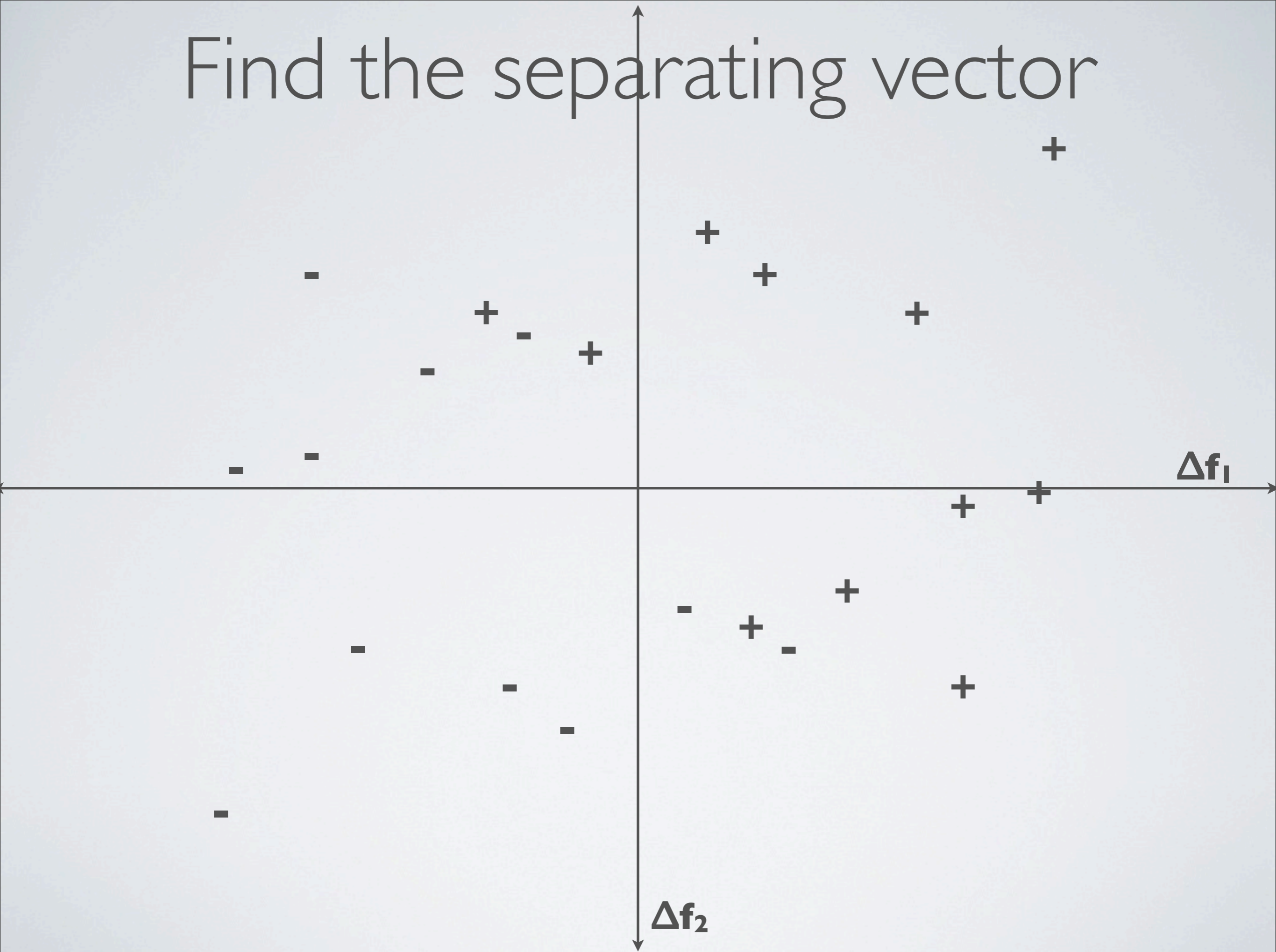
model

$$\vec{w} \cdot \vec{f_a - f_b} > 0$$

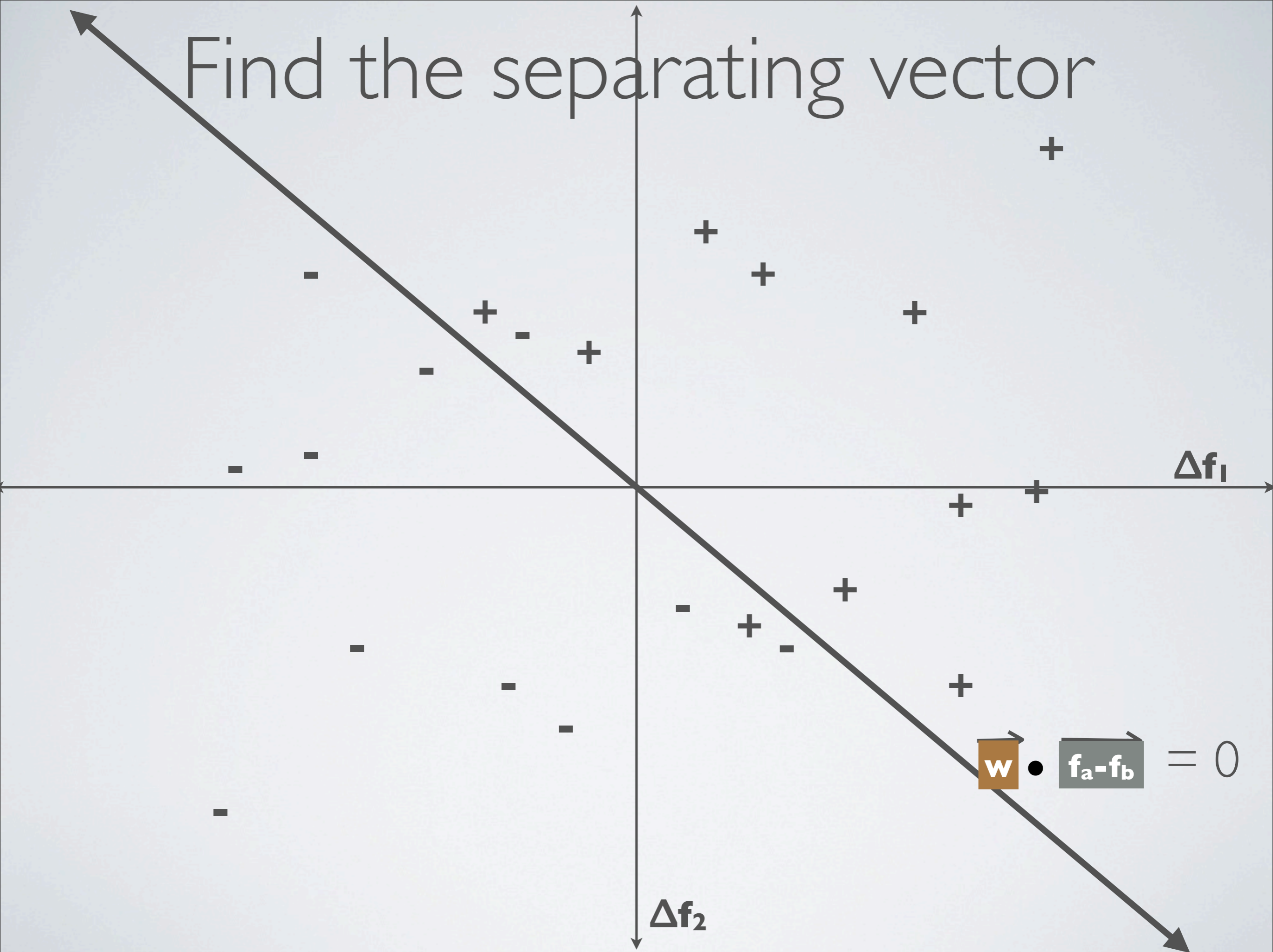
training instance
(**difference vector**)



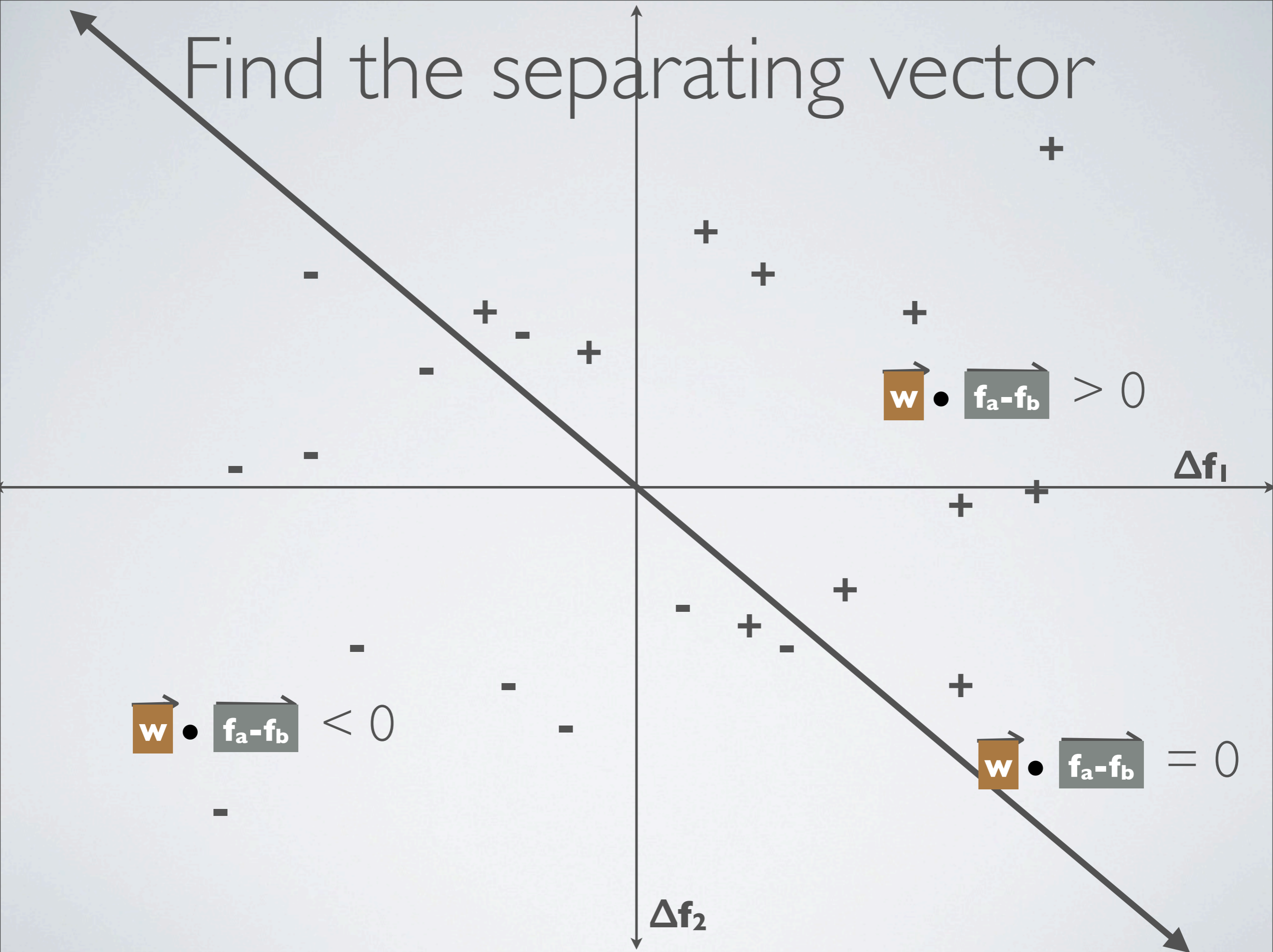
Find the separating vector



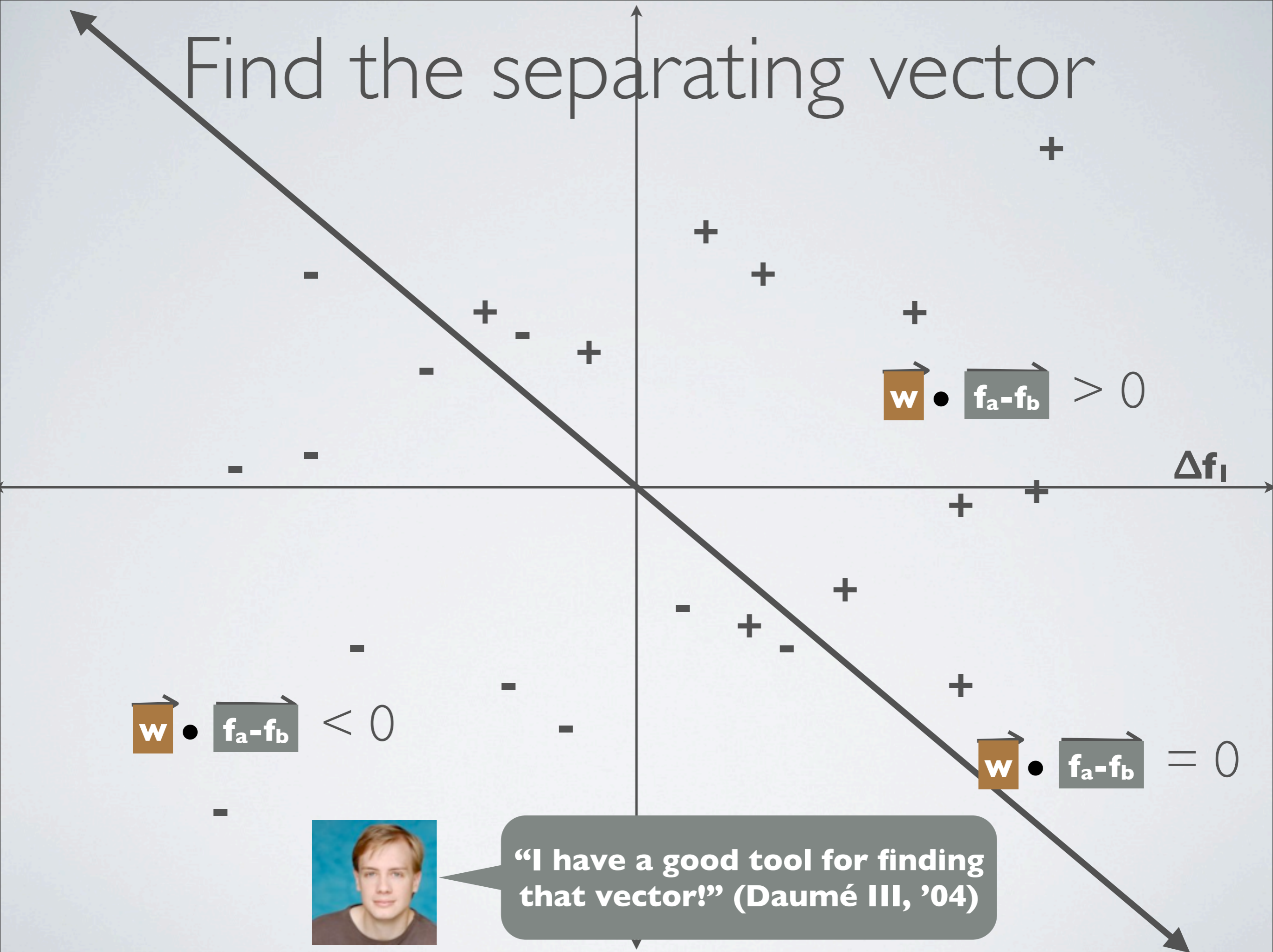
Find the separating vector



Find the separating vector



Find the separating vector



$$\vec{w} \cdot \vec{f_a - f_b} > 0$$

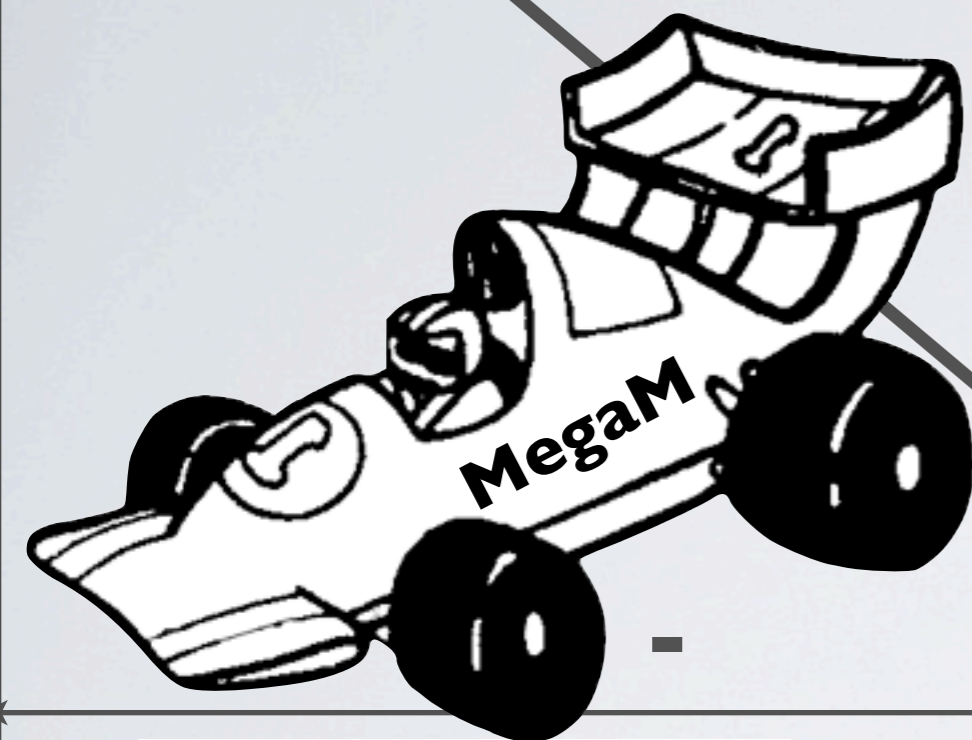
$$\vec{w} \cdot \vec{f_a - f_b} < 0$$

$$\vec{w} \cdot \vec{f_a - f_b} = 0$$



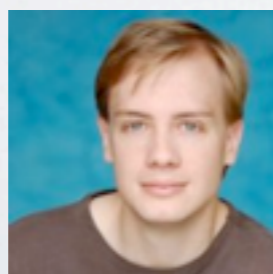
“I have a good tool for finding that vector!” (Daumé III, '04)

Find the separating vector



Daumé III, '04

$$\vec{w} \cdot \vec{f_a - f_b} < 0$$



"I have a good tool for finding that vector!" (Daumé III, '04)

$$\vec{w} \cdot \vec{f_a - f_b} > 0$$

$$\vec{w} \cdot \vec{f_a - f_b} = 0$$

Δf_1

Find the separating vector

+

+

$$\vec{w} \cdot \vec{f_a - f_b} > 0$$

Δf_1

+

+

Daumé III, '04

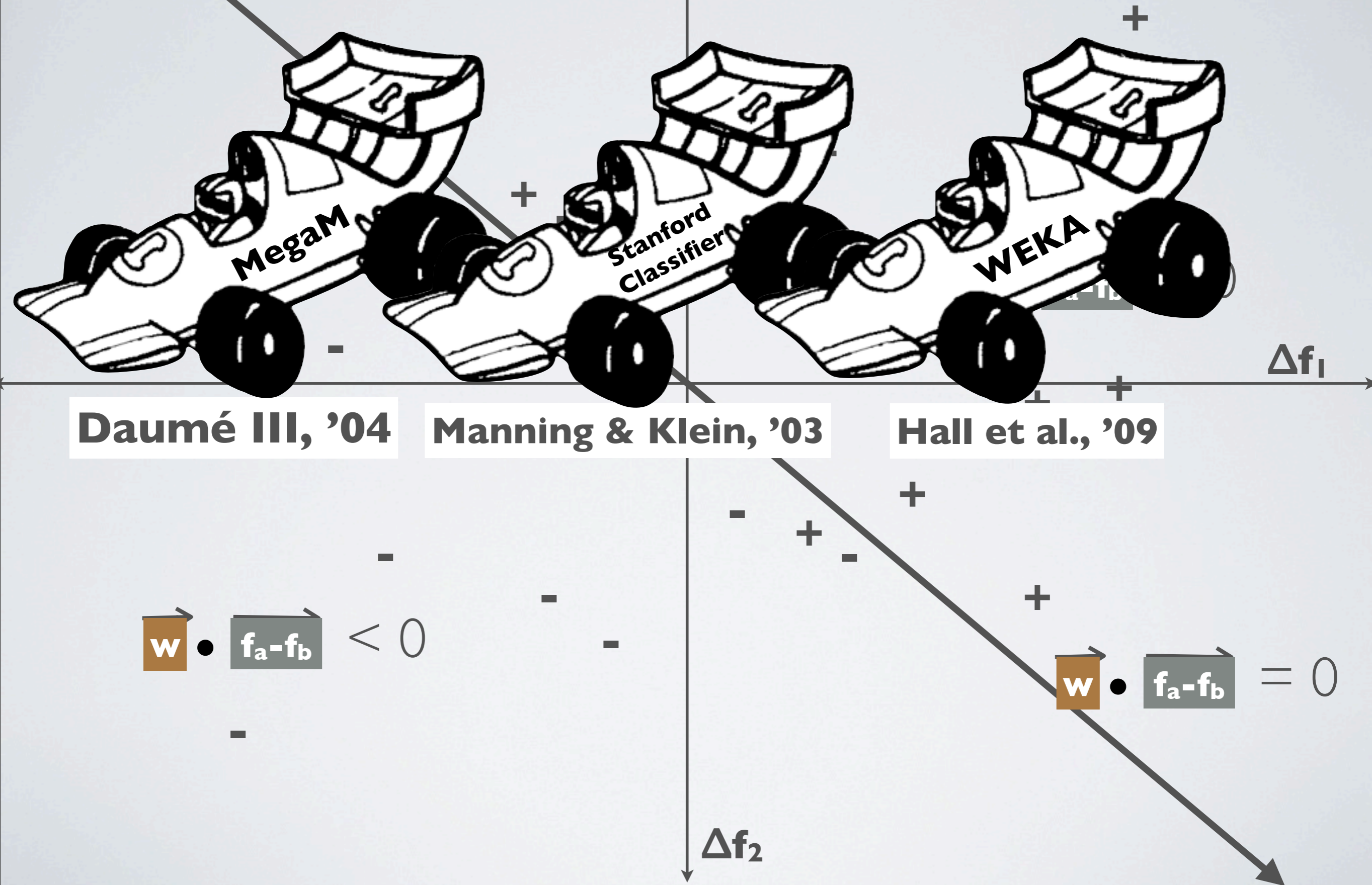
Manning & Klein, '03

$$\vec{w} \cdot \vec{f_a - f_b} < 0$$

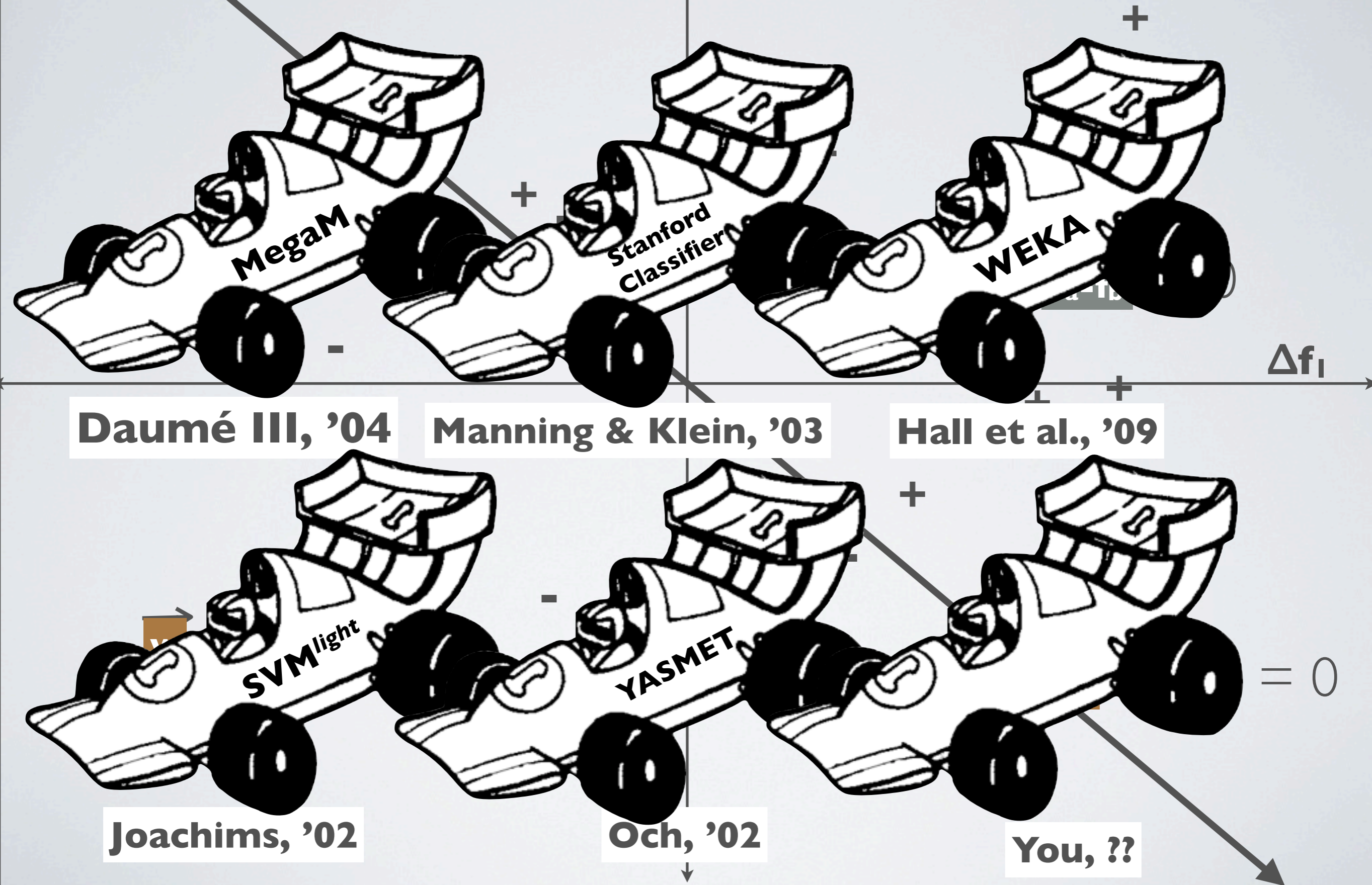
$$\vec{w} \cdot \vec{f_a - f_b} = 0$$

Δf_2

Find the separating vector



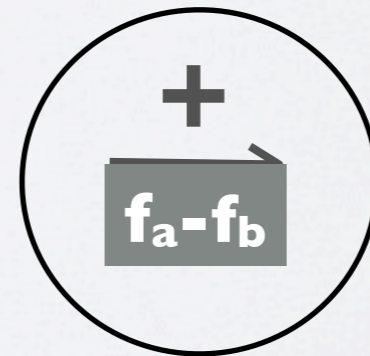
Find the separating vector



Avoid Intractability



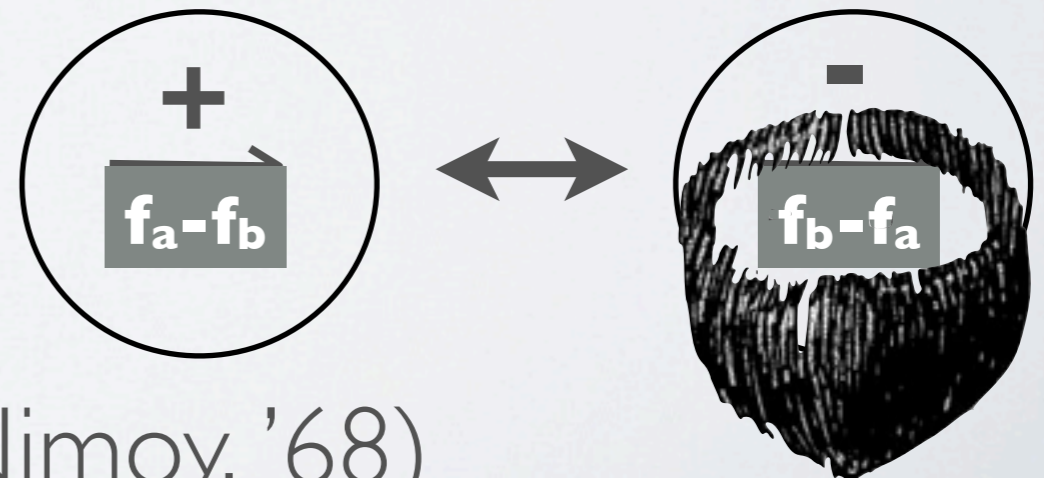
- Sample from the pool to avoid blowup
- Focus on difference vectors with large differences
- Add evil twins to ensure balance



Avoid Intractability

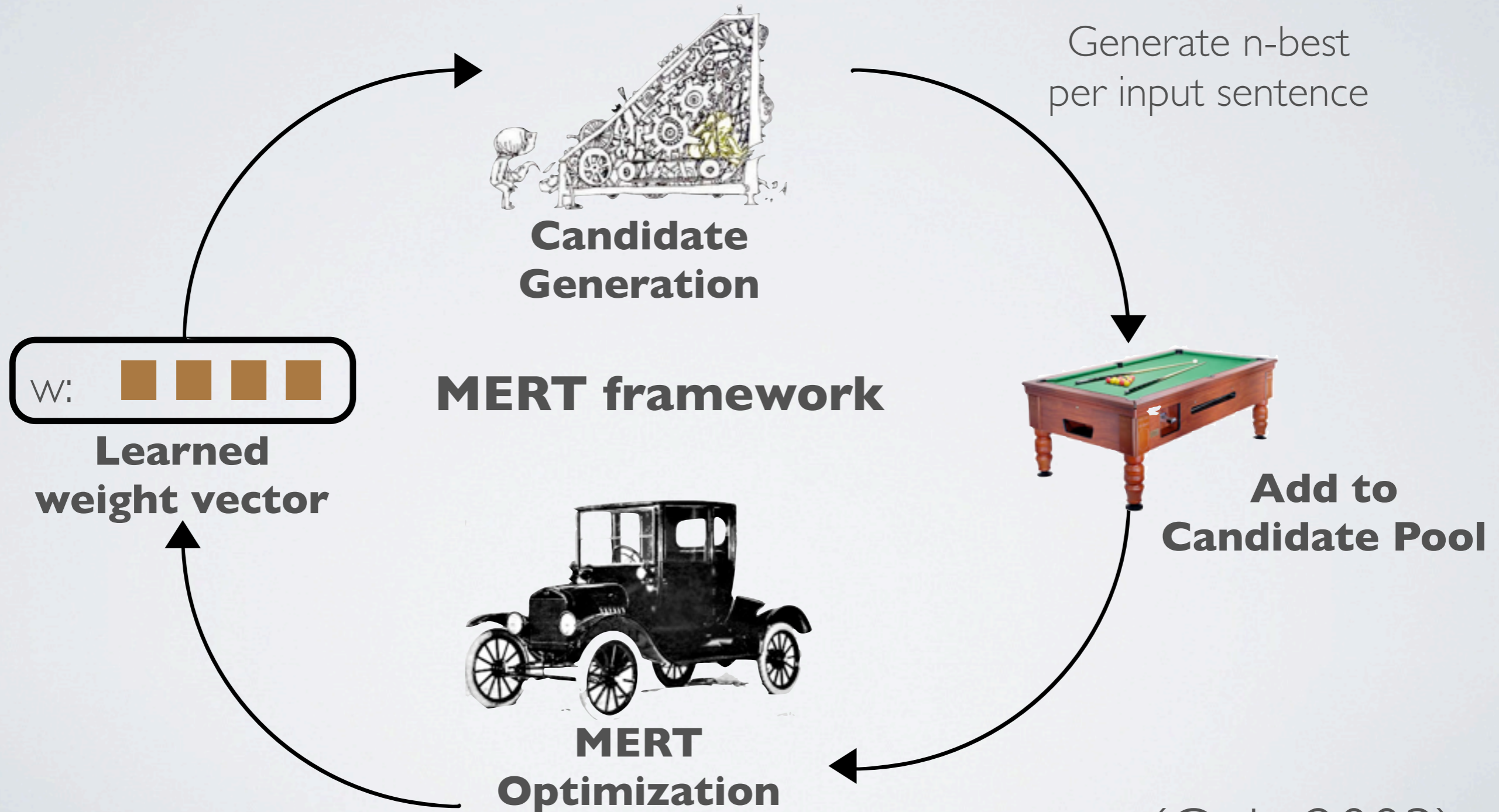


- Sample from the pool to avoid blowup
- Focus on difference vectors with large differences
- Add evil twins to ensure balance

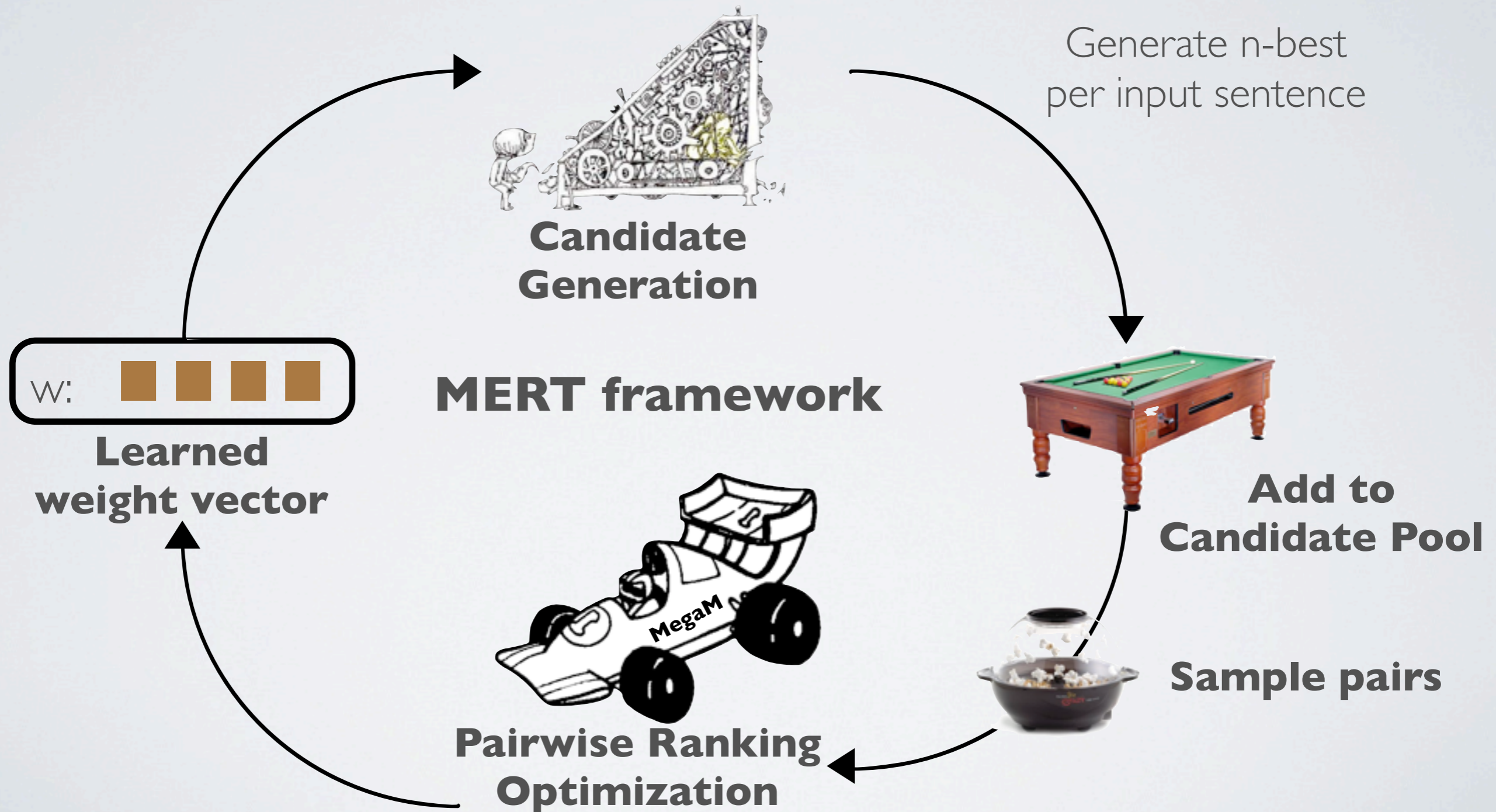


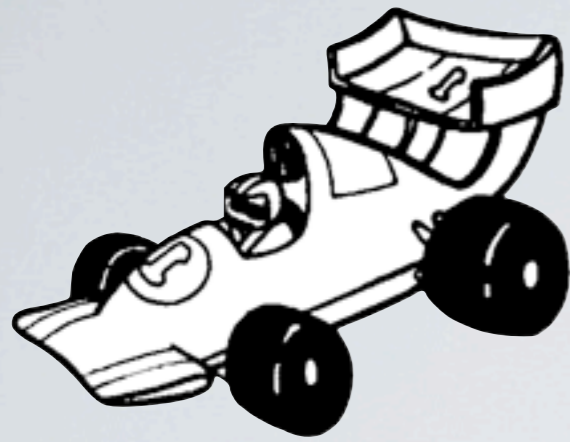
(Nimoy, '68)

MERT Tuning



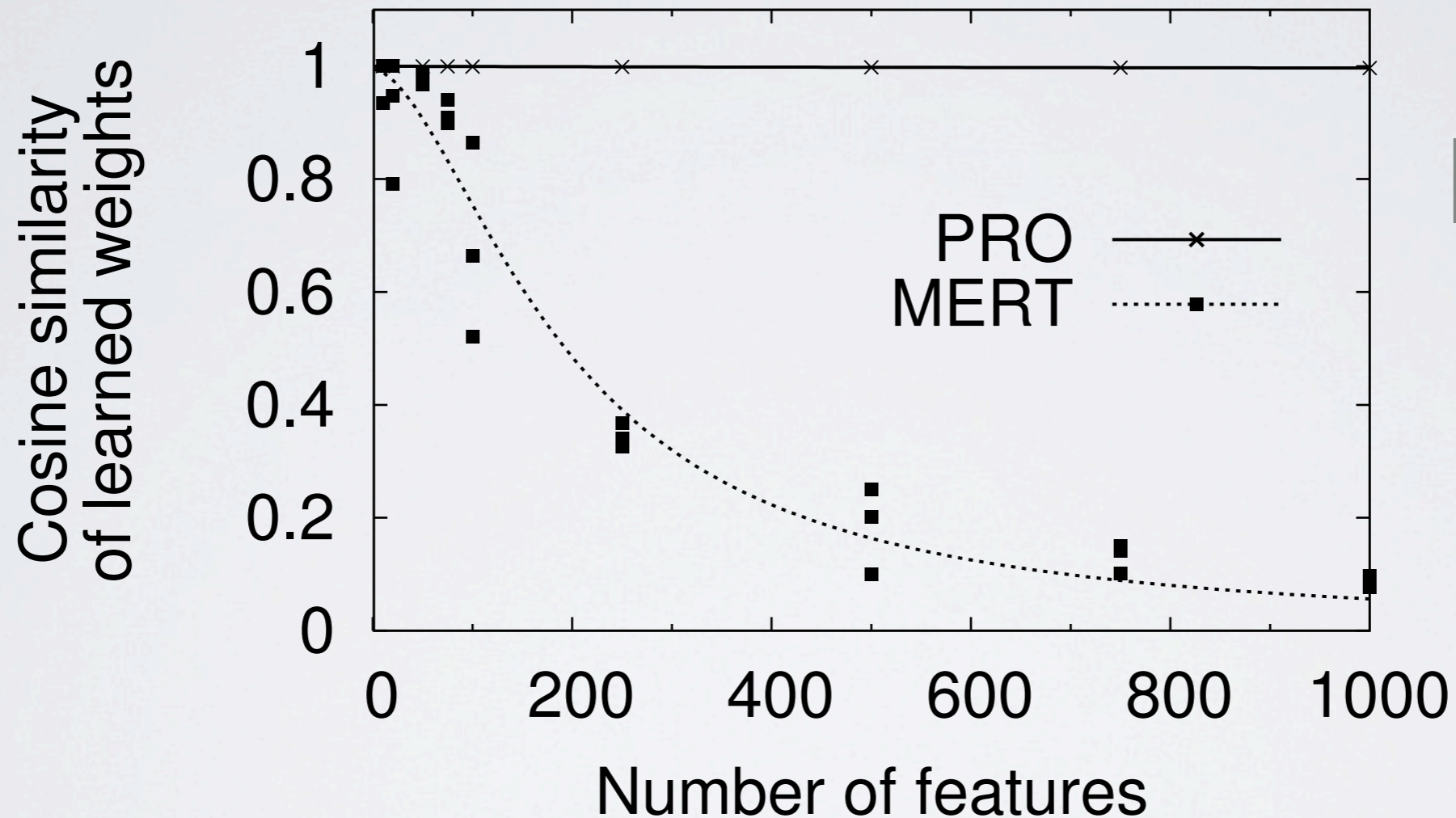
Pairwise Ranking Optimization (PRO) Tuning



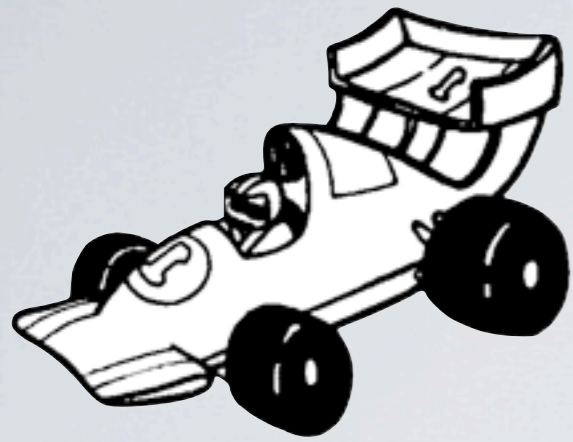


PRO *scales*

Synthetic weight learning
of MERT and PRO

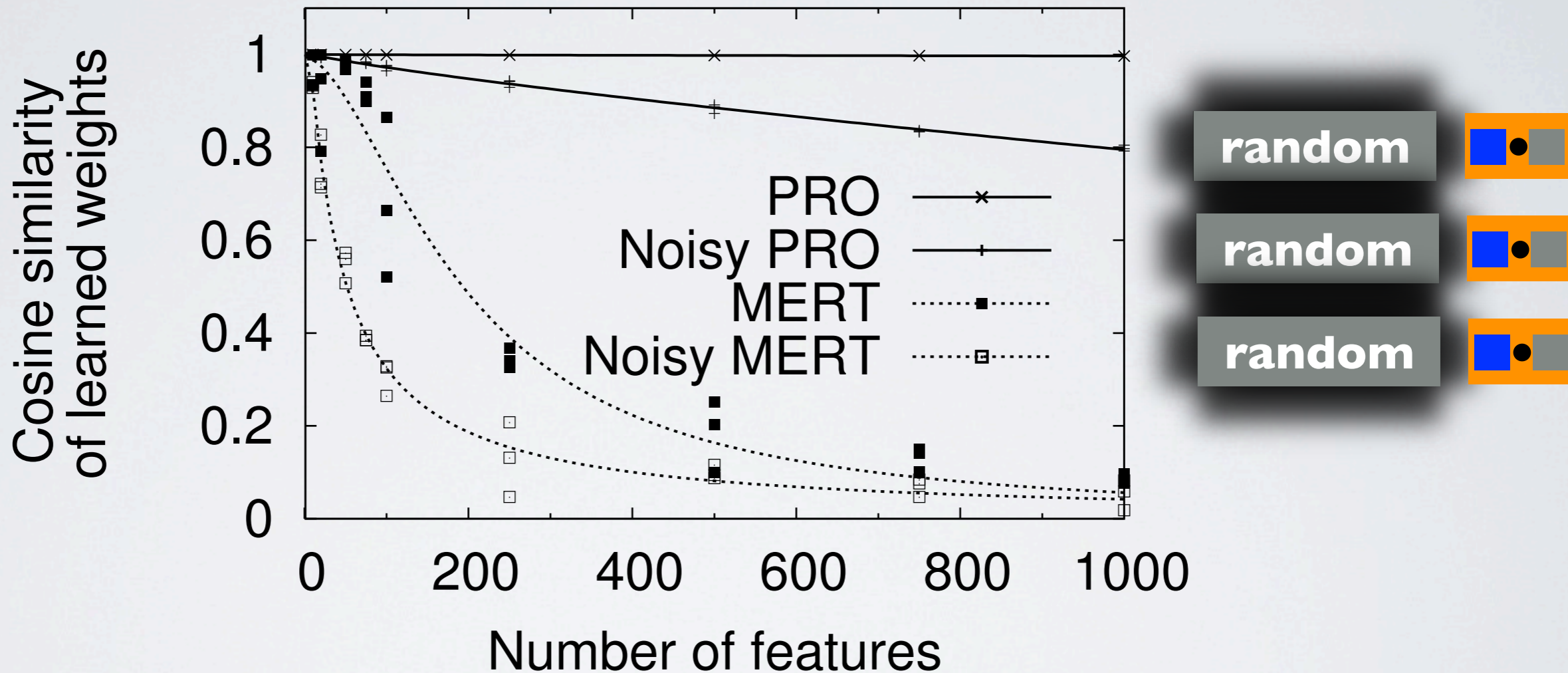


Unlike MERT, PRO is unfazed by a large number of features in the synthetic test



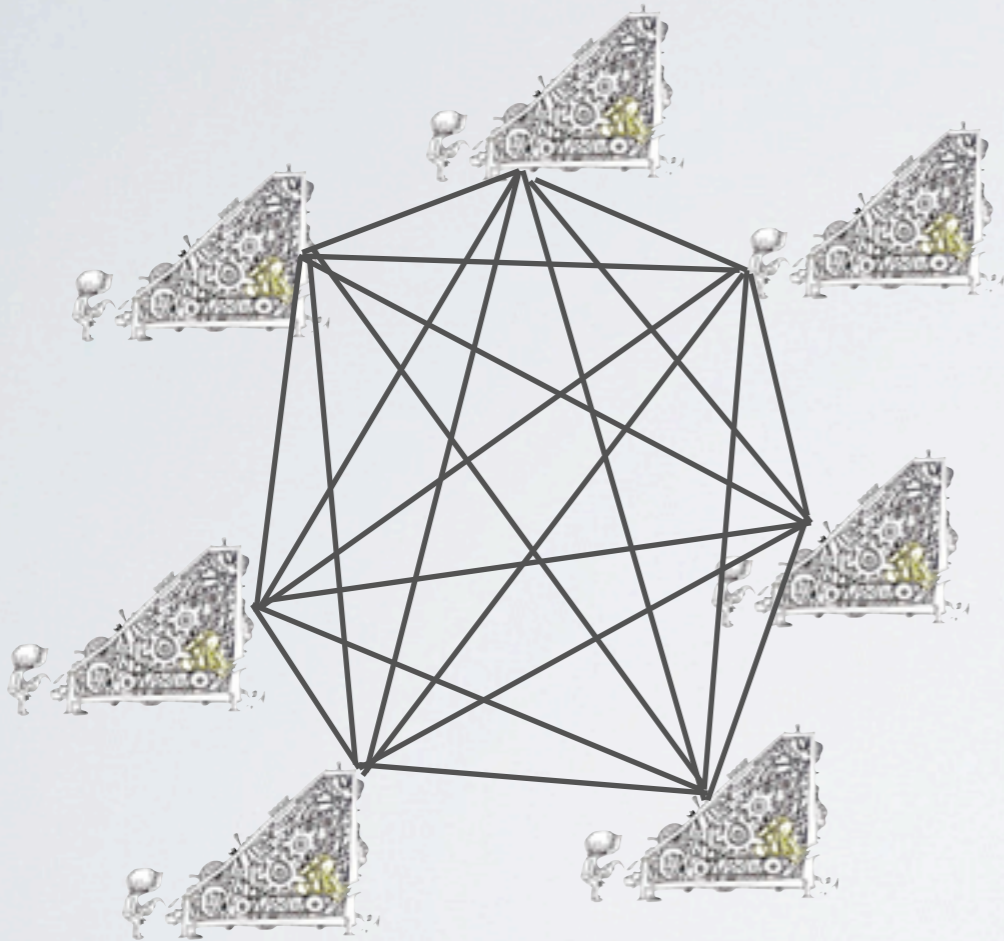
PRO *scales*

Synthetic weight learning
of MERT and PRO



Adding noise to the synthetic test makes it more difficult
but PRO still does quite well compared to MERT

MIRA also scales... but it's **hard** to implement



- Like PRO, a discriminative learning algorithm
- Unlike PRO, requires online, simultaneous optimization and decoding
- MIRA tuning must be customized to compute environment (cluster, inter-process communication, reliability concerns)

(Watanabe et al., '07)

(Chiang et al., '08, '09)

Unavoidable slide detailing the configuration and data of the experimental conditions....ZZZZZ

Language	Data (words)			Features			
	Train	Tune	Test	PBMT		SBMT	
				base	ext	base	ext
Urdu-English	2.2M (NIST 2009)	16K (NIST 2008)	18K (NIST 2008)	15	2250	19	277
Arabic-English	175M (NIST 2008)	65K (NIST 03-06/GALE)	47K (NIST 2008)	15	6333	19	352
Chinese-English	173M (GALE 2008)	42K (NIST 03-06)	37K (NIST 2008)	15	1828	19	517

Unavoidable slide detailing the configuration and data of the experimental conditions....ZZZZZ

Language	Data (words)		Evaluation				
	Train	Tune	base	ext	base	ext	
Urdu-English	2.2M (NIST 2009)	16K (NIST 2008)	18K (NIST 2008)	15	2250	19	277
Arabic-English	175M (NIST 2008)	65K (NIST 03-06/GALE)	47K (NIST 2008)	15	6333	19	352
Chinese-English	173M (GALE 2008)	42K (NIST 03-06)	37K (NIST 2008)	15	1828	19	517

Standard large-scale language pairs

2.2M
(NIST 2009)

175M
(NIST 2008)

173M
(GALE 2008)

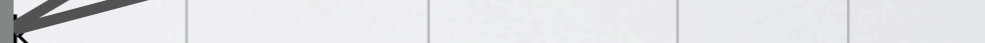
Unavoidable slide detailing the configuration and data of the experimental conditions....ZZZZZ

Language	Data (words)			Features			
	Train	Tune	Test	PBMT		SBMT	
				base	ext	base	ext
Urdu-English	2.1M (NIST 2008)	65K (NIST 2008)	47K (NIST 2008)	15	2250	19	277
Arabic-English	175M (NIST 2008)	65K (NIST 03-06/GALE)	47K (NIST 2008)	15	6333	19	352
Chinese-English	173M (GALE 2008)	42K (NIST 03-06)	37K (NIST 2008)	15	1828	19	517

State-of-the-art decoders

PBMT

SBMT



Unavoidable slide detailing the configuration and data of the experimental conditions...ZZZZZ

Two feature configurations per decoder

Language	Train	Tune	Test	Features			
				PBMT		SBMT	
				base	ext	base	ext
Urdu-English	2.2M (NIST 2009)	16K (NIST 2008)	18K (NIST 2008)	15	2250	19	277
Arabic-English	175M (NIST 2008)	65K (NIST 03-06/GALE)	47K (NIST 2008)	15	6333	19	352
Chinese-English	173M (GALE 2008)	42K (NIST 03-06)	37K (NIST 2008)	15	1828	19	517

Unavoidable slide detailing the configuration and data of the experimental conditions....ZZZZZ

Language	Data (words)			Features			
	Ran MERT, MIRA, PRO			PBMT		SBMT	
				base	ext	base	ext
Urdu-English	2.2M (NIST 2009)	16K (NIST 2008)	18K (NIST 2008)	15	2250	19	277
Arabic-English	175M (NIST 2008)	65K (NIST 03-06/GALE)	47K (NIST 2008)	15	6333	19	352
Chinese-English	173M (GALE 2008)	42K (NIST 03-06)	37K (NIST 2008)	15	1828	19	517

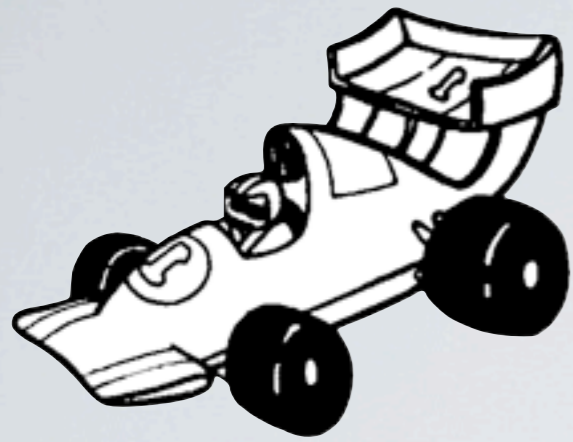
Unavoidable slide detailing the configuration and data of the experimental conditions....zzzzz

Language	Data (words)			Features			
	<div style="background-color: #cccccc; padding: 10px; text-align: center;"> Ran MIRA, PRO (MERT doesn't scale) </div>			PBMT		SBMT	
				base	ext	base	ext
Urdu-English	2.2M (NIST 2009)	16K (NIST 2008)	18K (NIST 2008)	15	2250	19	277
Arabic-English	175M (NIST 2008)	65K (NIST 03-06/GALE)	47K (NIST 2008)	15	6333	19	352
Chinese-English	173M (GALE 2008)	42K (NIST 03-06)	37K (NIST 2008)	15	1828	19	517

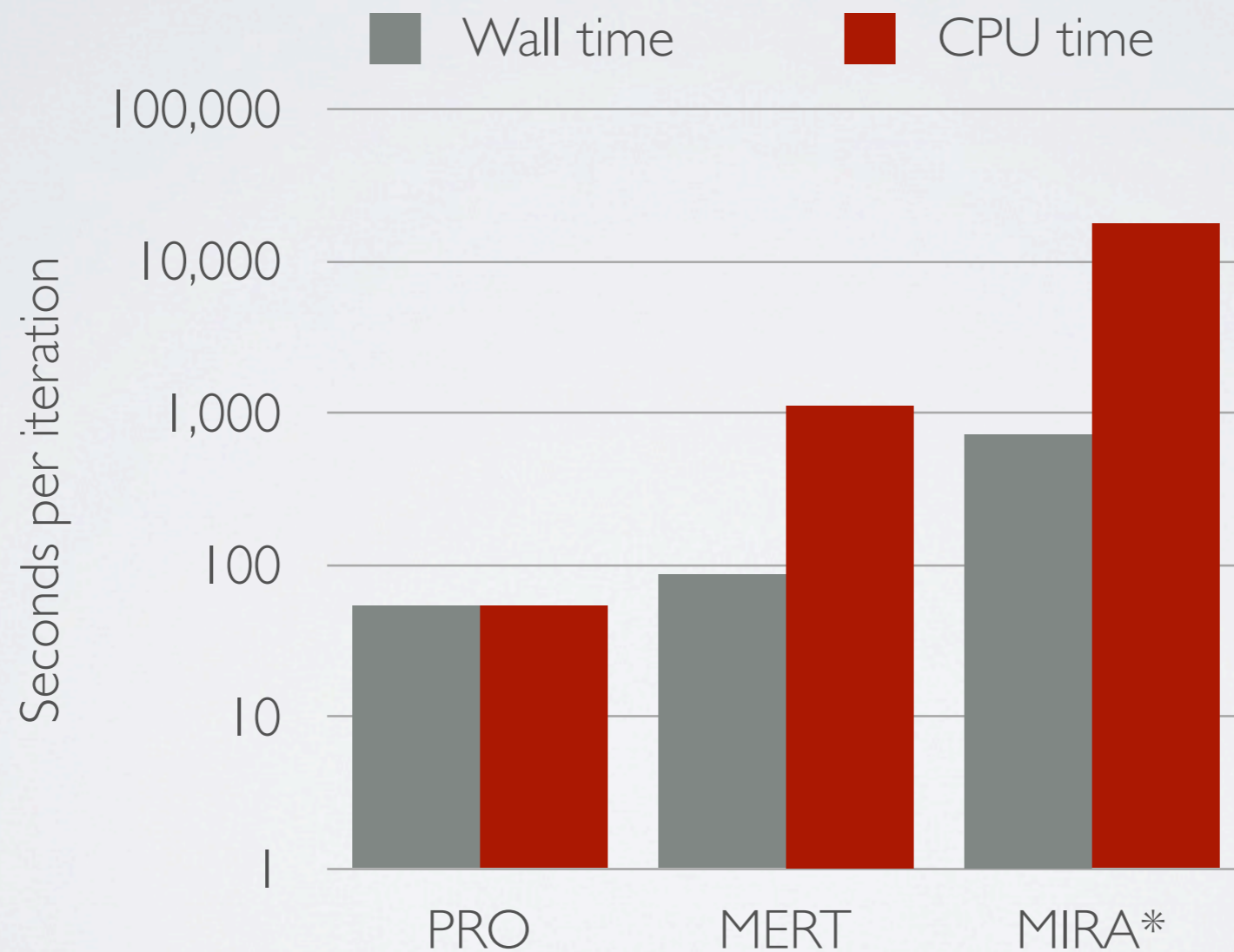
Unavoidable slide detailing the configuration experimental

Report 4-reference, detokenized, mixed-case BLEU

Language	Data (words)			Features			
	Train	Tune	Test	PBMT		SBMT	
				base	ext	base	ext
Urdu-English	2.2M (NIST 2009)	16K (NIST 2008)	18K (NIST 2008)	15	2250	19	277
Arabic-English	175M (NIST 2008)	65K (NIST 03-06/GALE)	47K (NIST 2008)	15	6333	19	352
Chinese-English	173M (GALE 2008)	42K (NIST 03-06)	37K (NIST 2008)	15	1828	19	517



PRO is *fast*

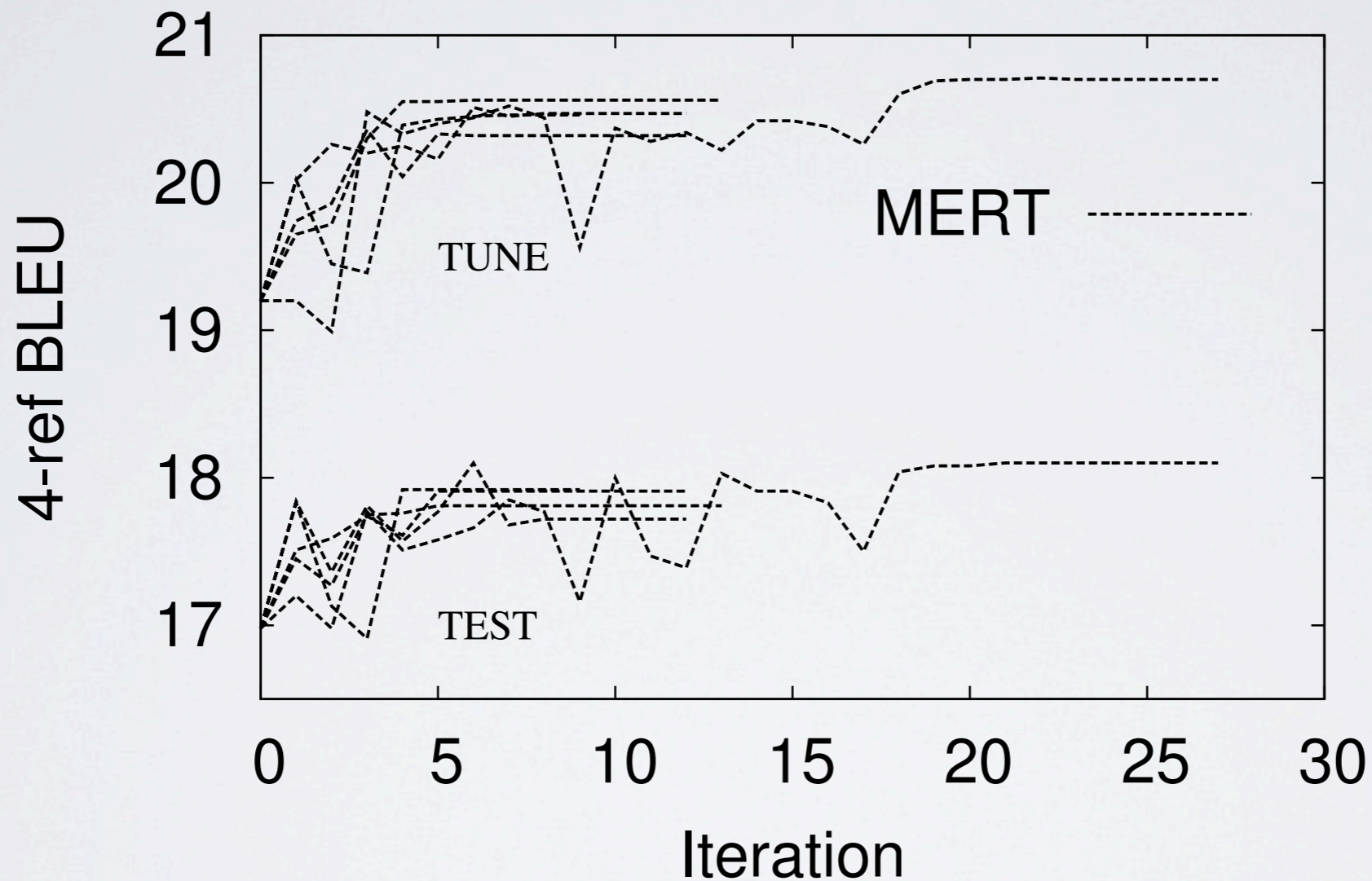


* Your implementation of MIRA may be faster



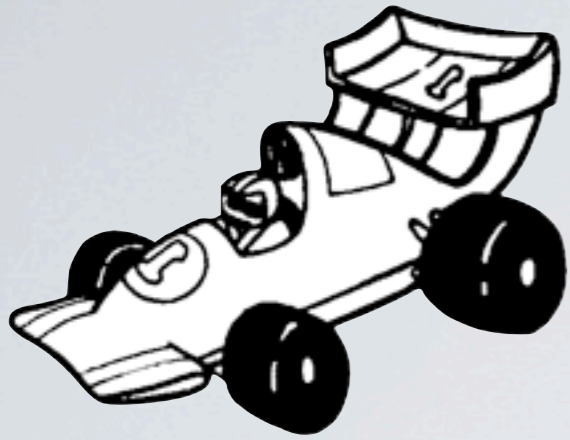
MERT is *unstable*

Urdu-English PBMT tuning stability



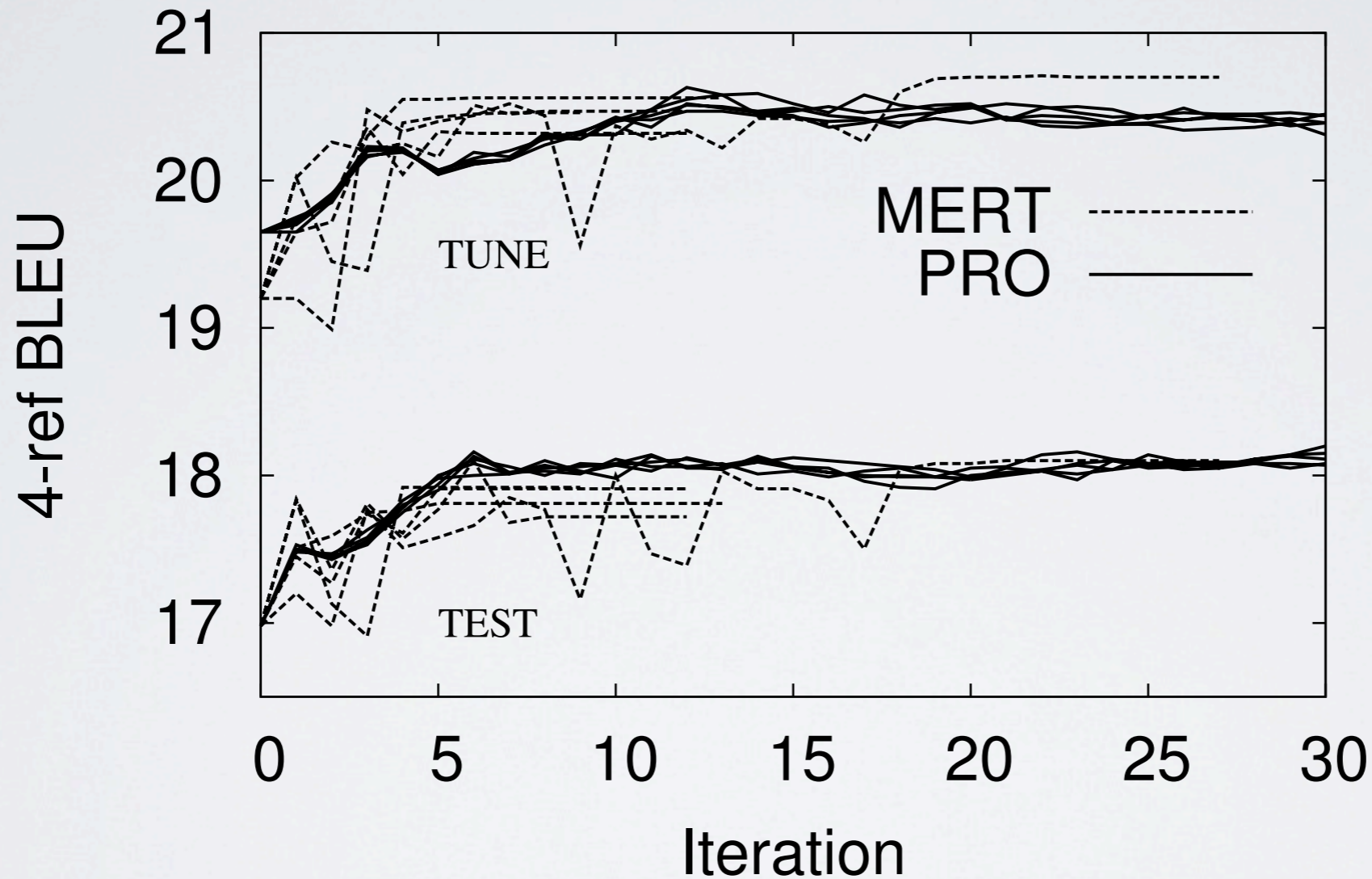
Result from five identical runs

(Clark et al., 2011)



PRO is *stable*

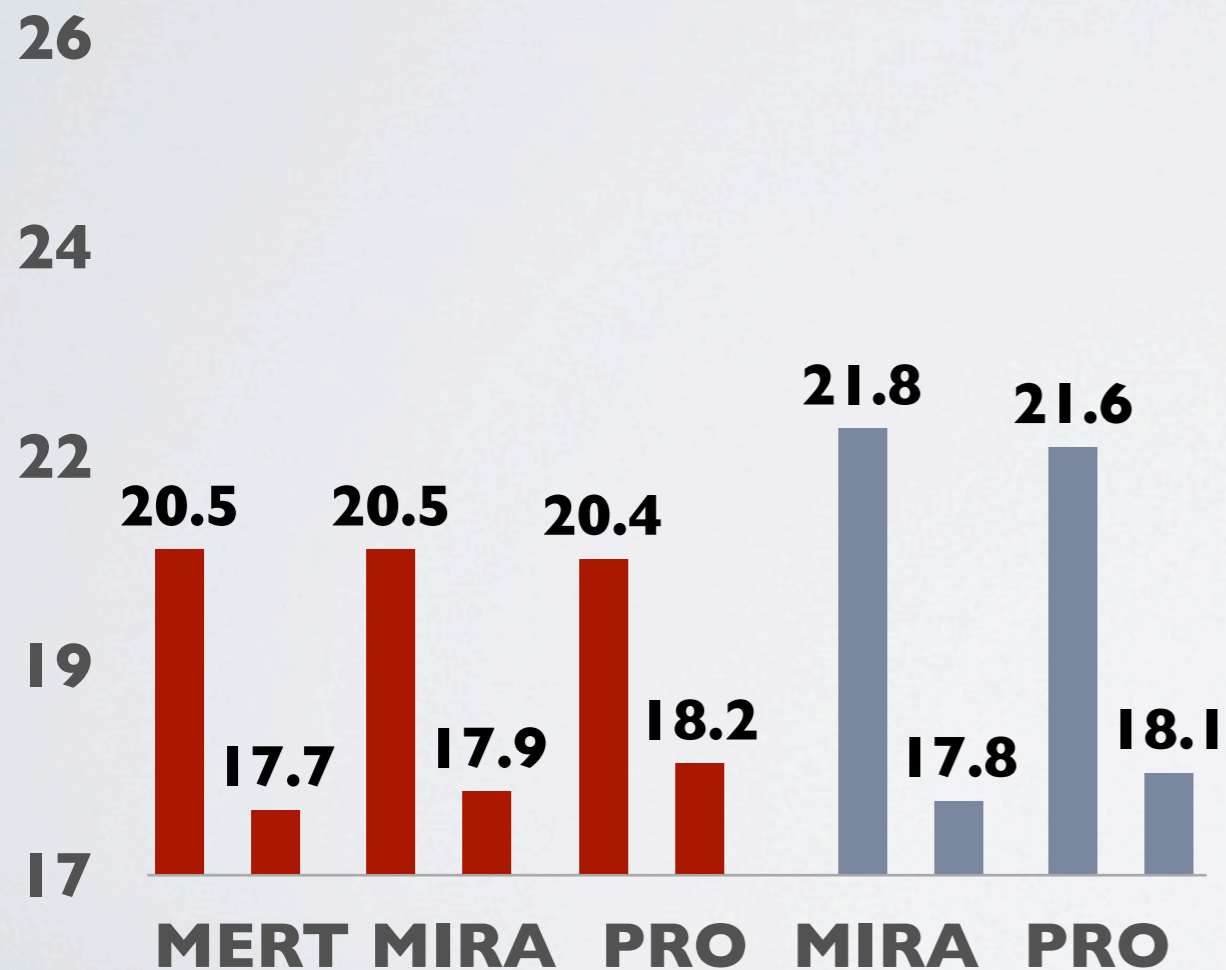
Urdu-English PBMT tuning stability



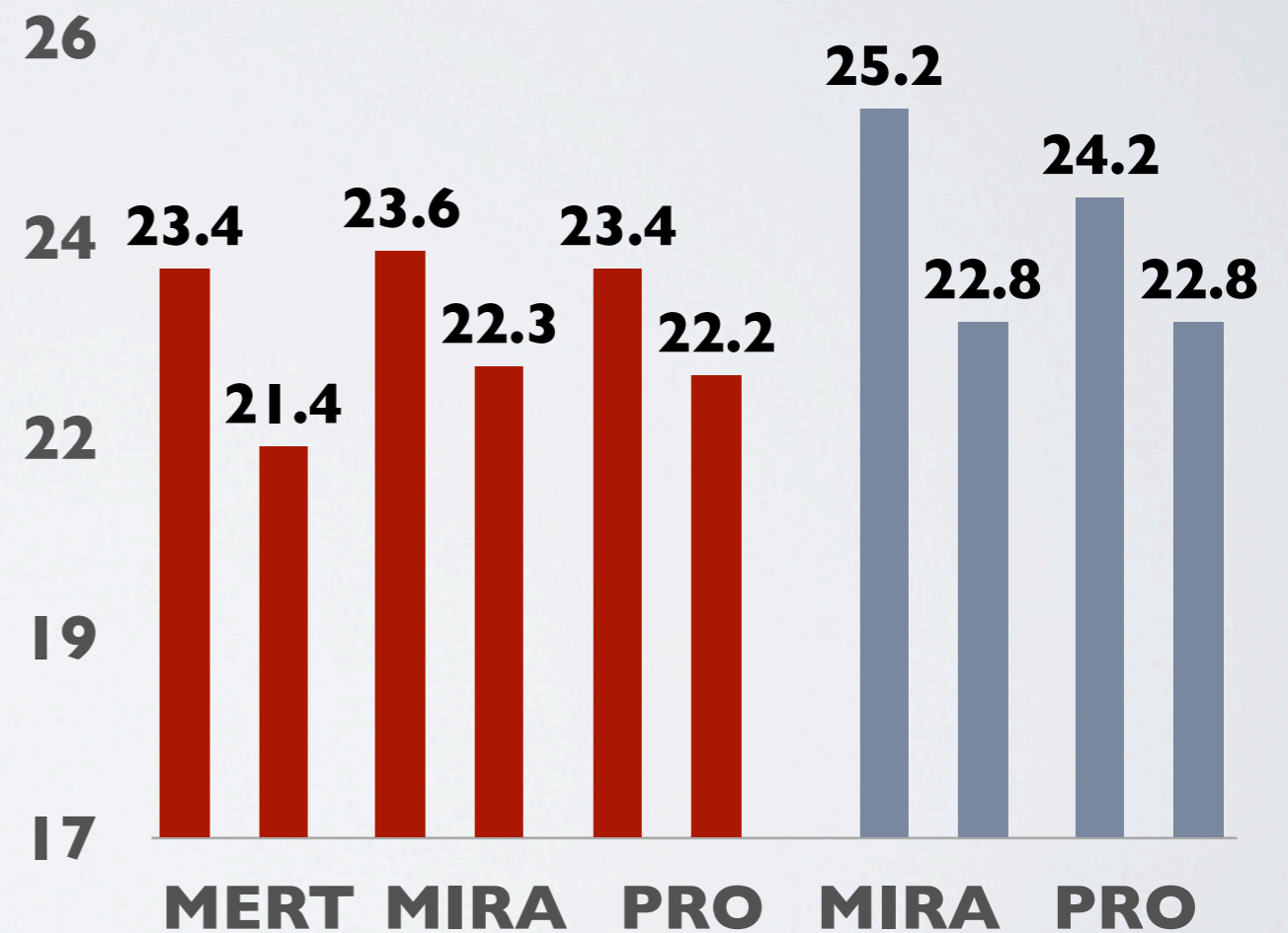
Result from five identical runs

MERT vs. MIRA vs. *PRO*

PBMT Urdu-English

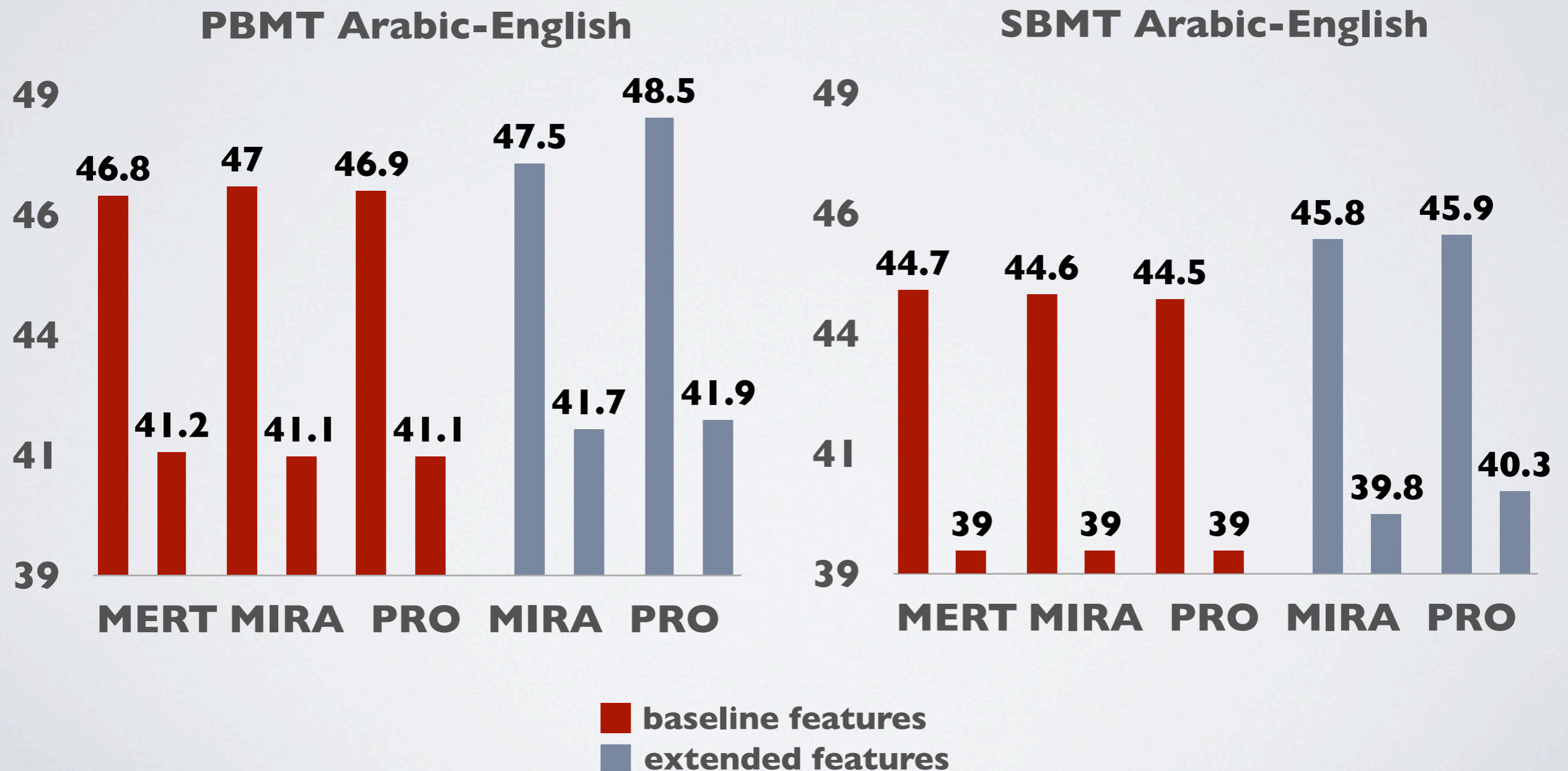


SBMT Urdu-English



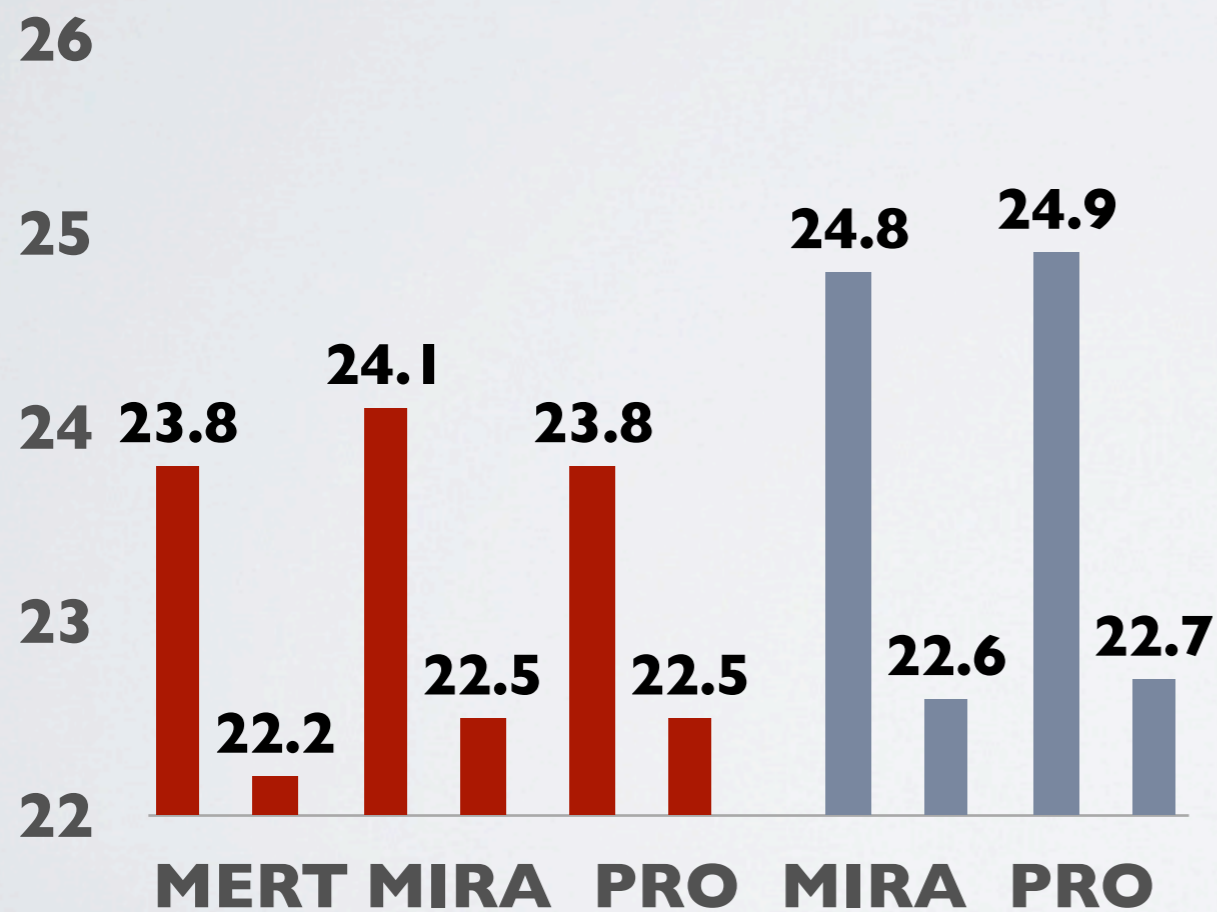
■ baseline features
■ extended features

MERT vs. MIRA vs. *PRO*

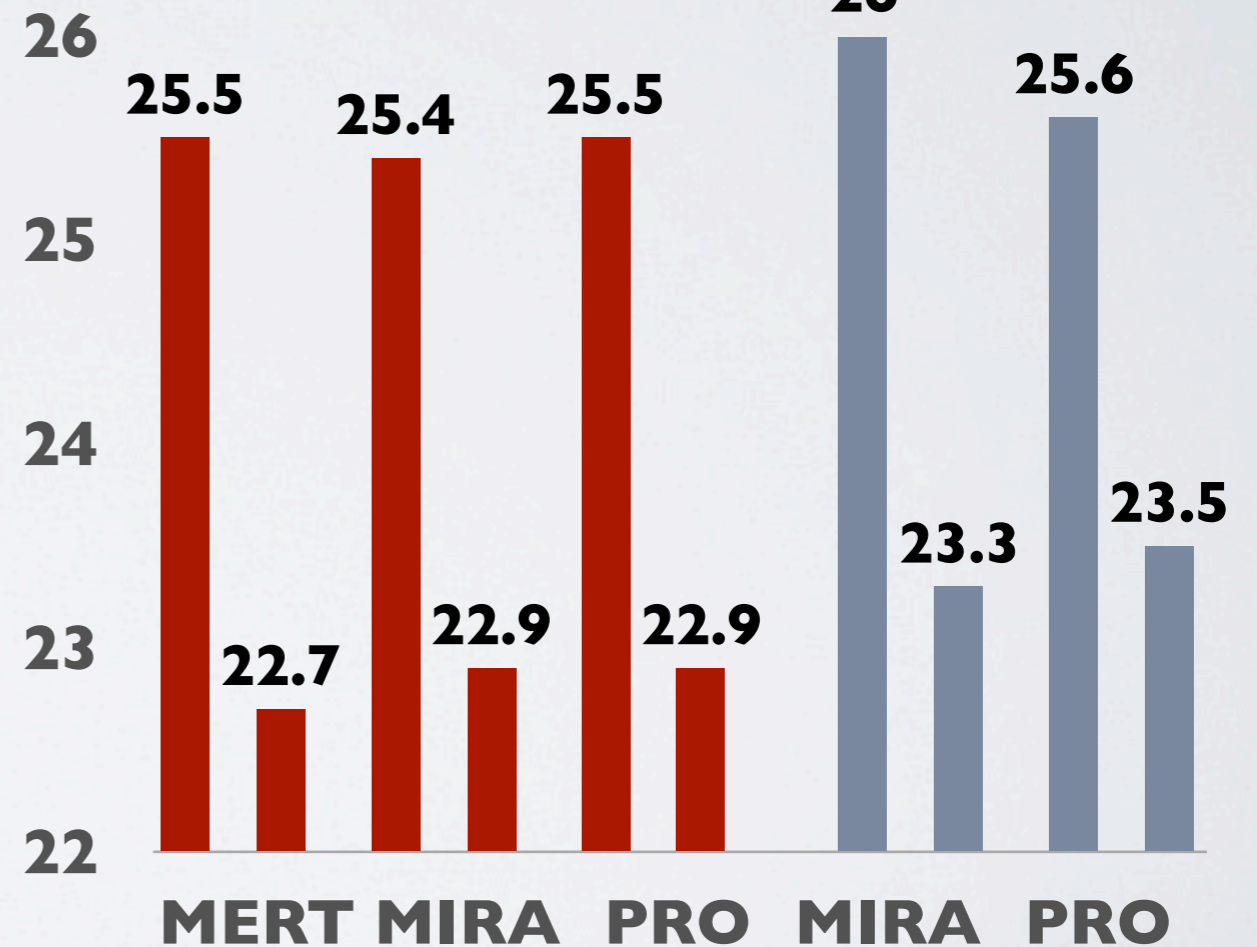


MERT vs. MIRA vs. *PRO*

PBMT Chinese-English



SBMT Chinese-English



■ baseline features
■ extended features

PRO is comparable to all

■ base
■ ext

PBMT Urdu-English



SBMT Urdu-English



PBMT Arabic-English



SBMT Arabic-English



PBMT Chinese-English



SBMT Chinese-English



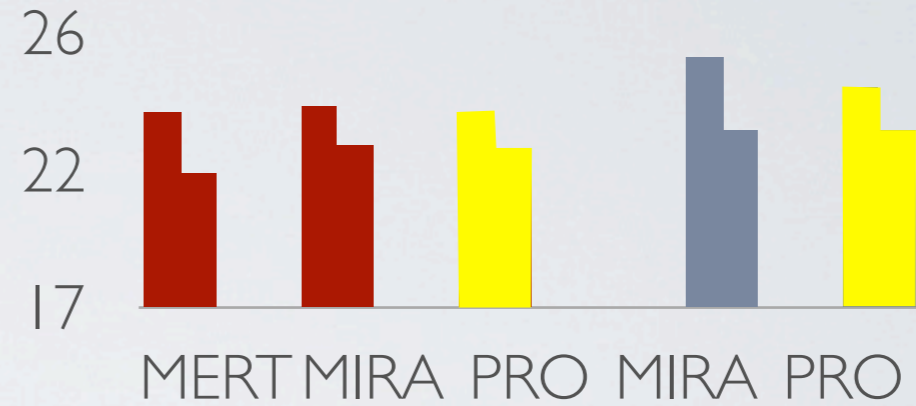
PRO is comparable to all

■ base
■ ext

PBMT Urdu-English



SBMT Urdu-English



PBMT Arabic-English



SBMT Arabic-English



PBMT Chinese-English



SBMT Chinese-English



Related Work

SampleRank

(Culotta, '08, Wick et al., '09, Roth et al., '10)

Similar approach, with guided search through pool space
(See Haddow et al. in WMT)

Classifier-based Weight Learning

(Tillmann & Zhang, '05, Och & Ney, '02
Ittycheriah & Roukos, '05, Xiong et al., '06)

Various approaches using classifiers to learn MT feature weights -- these do not use the difference vector approach

Discriminative Re-ranking

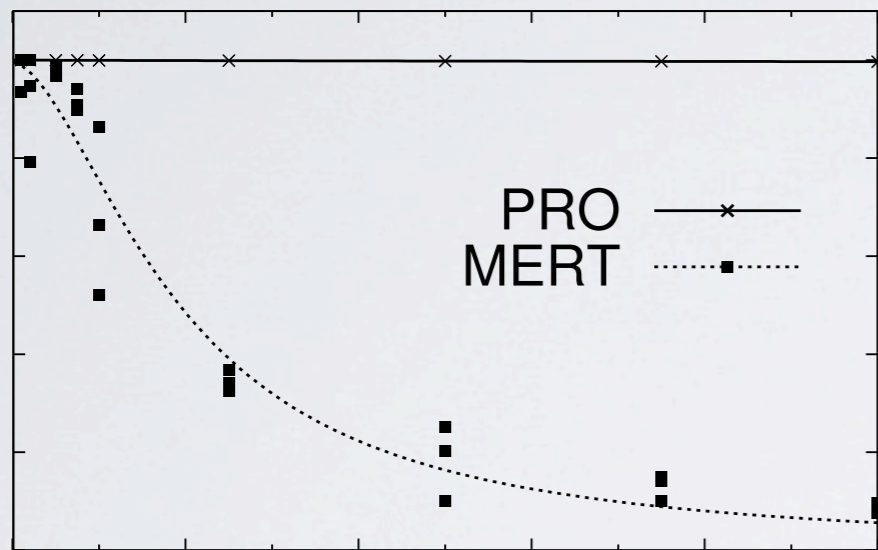
(Shen et al., '04, Cowan et al., '06,
Watanabe et al., '06)

Changing the n-best list after decoding using similar techniques to ours

Why Use **PRO**?

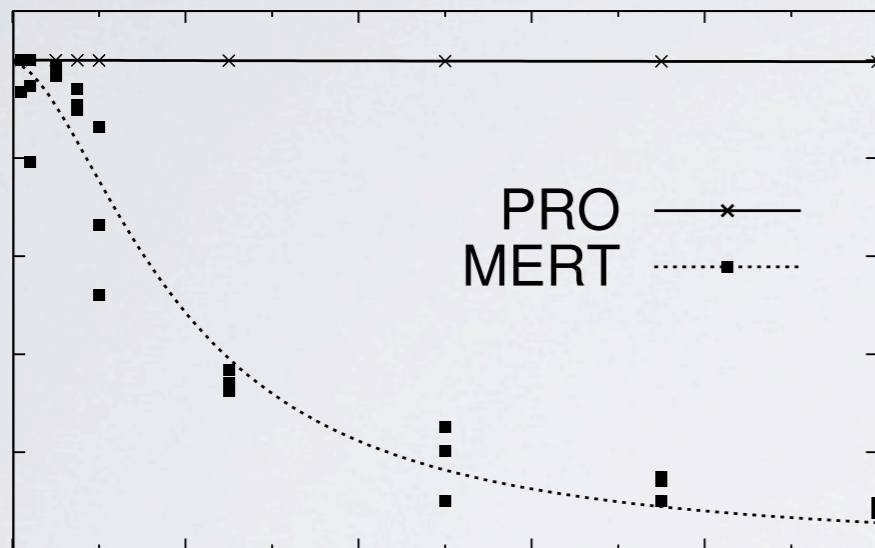
Why Use **PRO**?

It's **scalable**

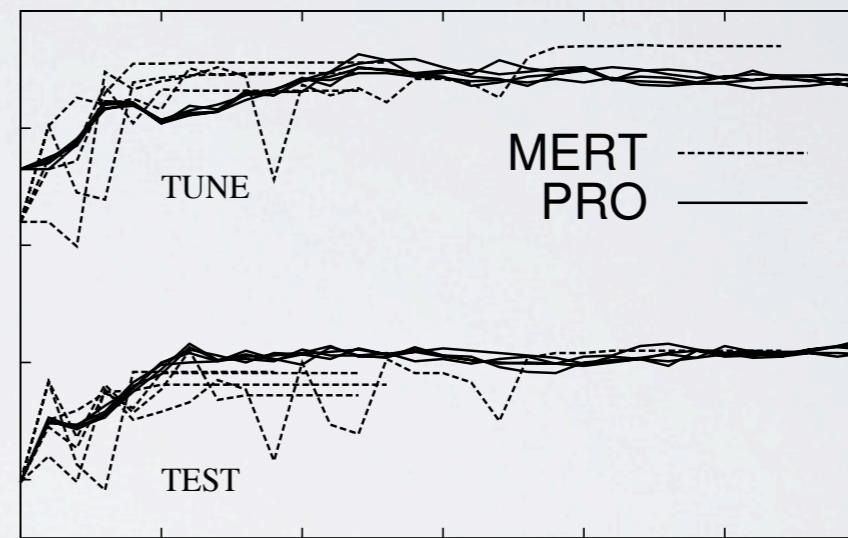


Why Use **PRO**?

It's **scalable**

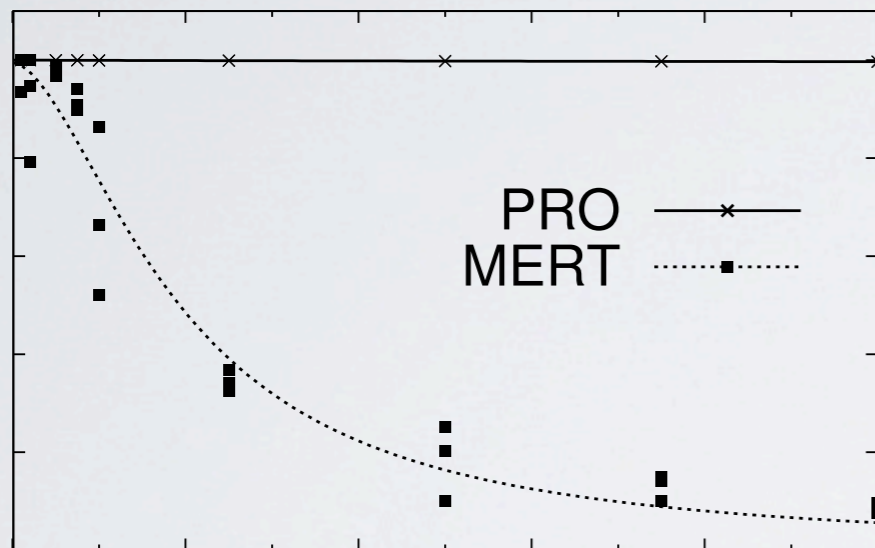


It's **stable**

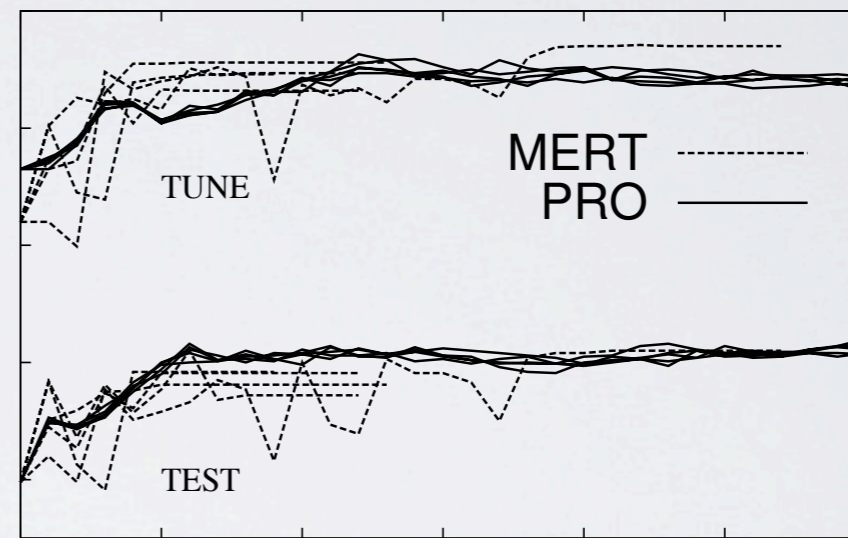


Why Use **PRO**?

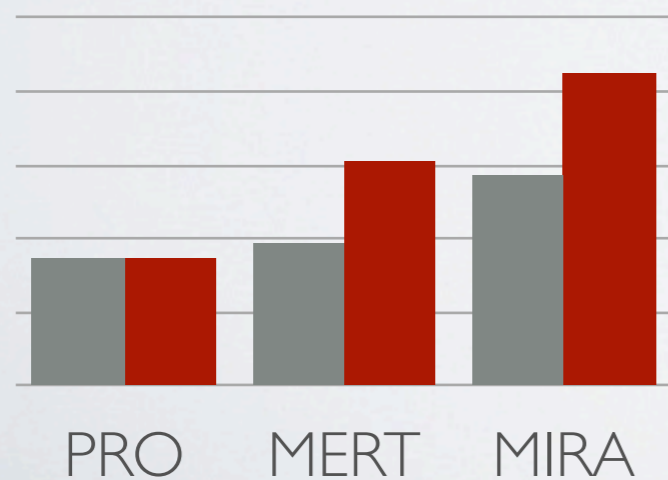
It's **scalable**



It's **stable**

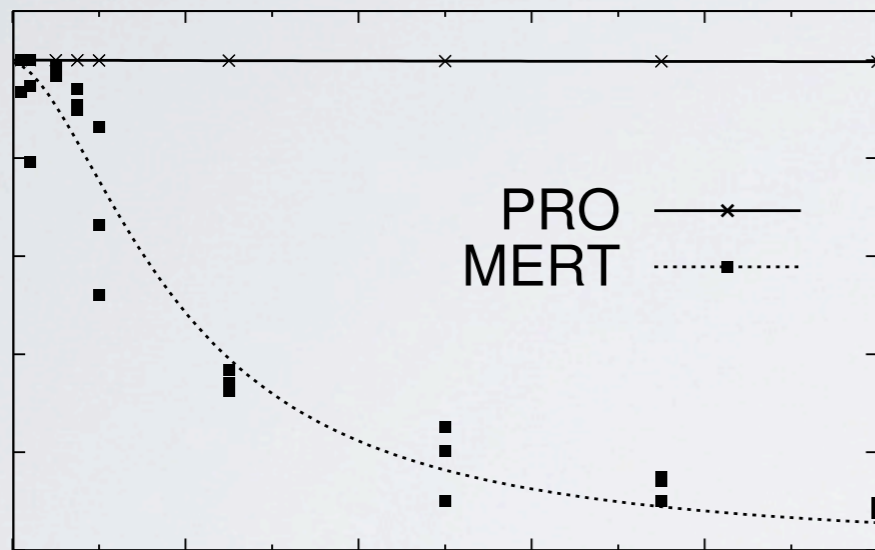


It's **fast**

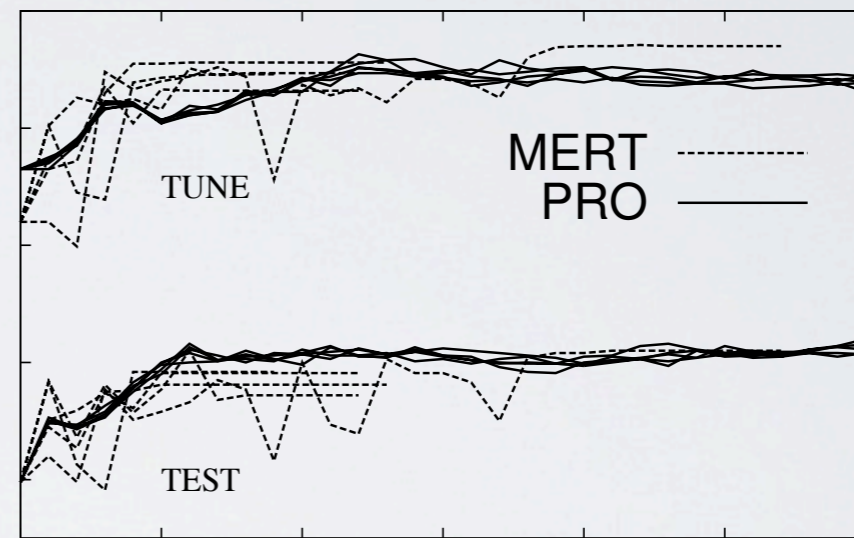


Why Use **PRO**?

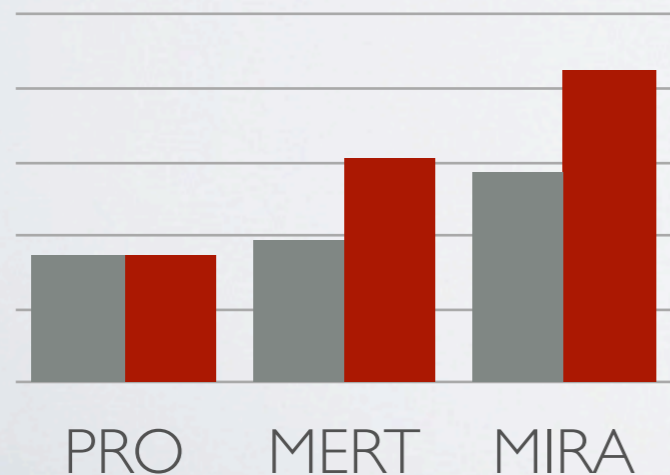
It's **scalable**



It's **stable**



It's **fast**

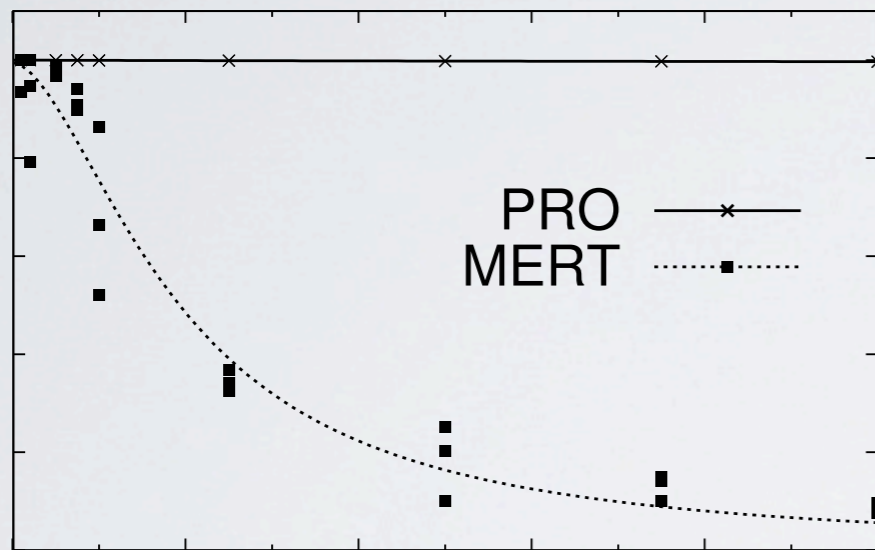


It's **easy**

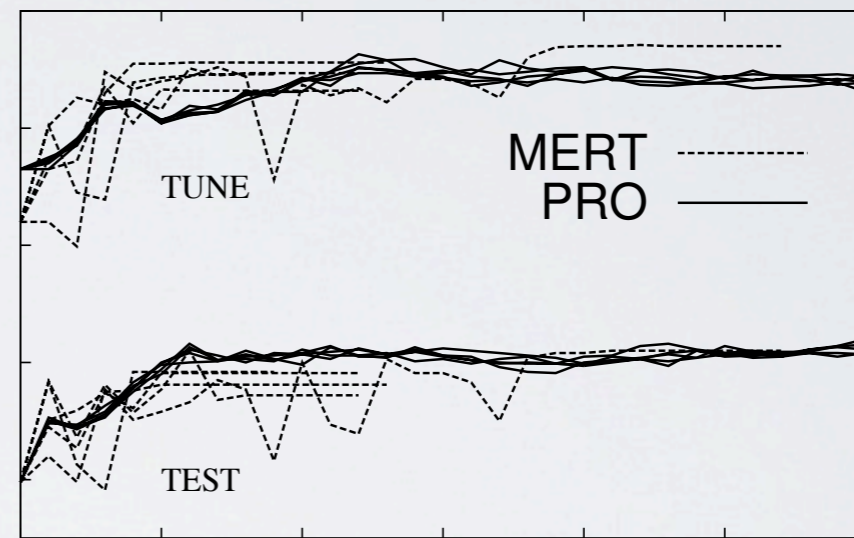
At least **three** external implementations prior to this talk

Why Use **PRO**?

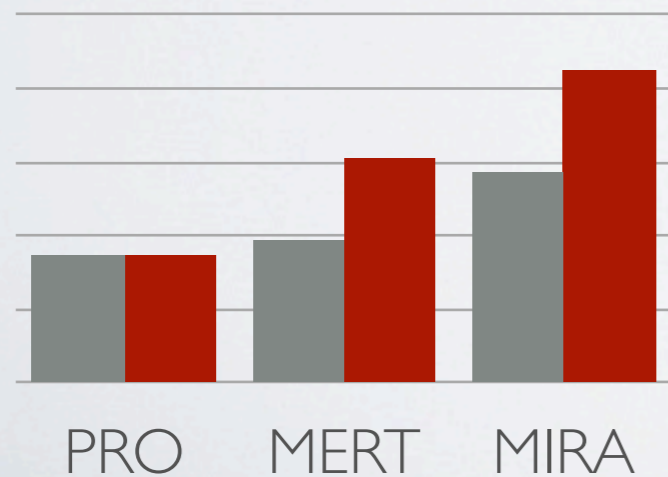
It's **scalable**



It's **stable**



It's **fast**



It's **easy**



**“Including mine!”
(Dyer, P.C.)**

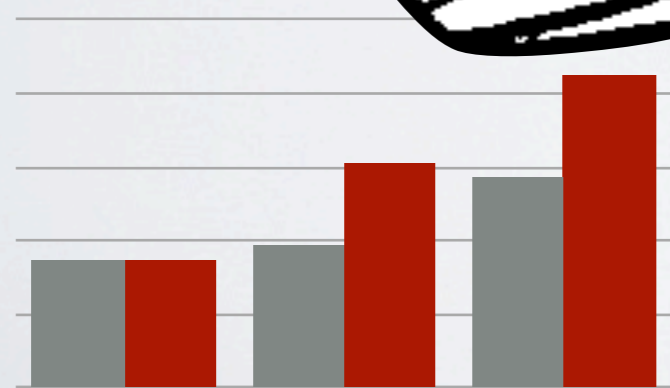
<https://github.com/redpony/cdec/tree/master/pro-train>

Why Use **PRO**?

It's **scalable**



It's **fast**



PRO

MERT

MIRA

It's **easy**



**“Including mine!”
(Dyer, P.C.)**

<https://github.com/redpony/cdec/tree/master/pro-train>