

第21回国際言語学オリンピック

ブラジル・ブラジリア、2024年7月23～31日

団体戦 問題

語彙統計学とは、語彙の類似性に基づいて言語間の近縁性を予測するために開発された一連の手法のことを指す。これらの手法は、通常、専門家によって手動でアノテーションを付けられた単語の長いリストに対して用いられ、専門家は特定の単語のペアが同じ語源に由来すると考えられるかどうかを表示する。しかしながら、たまに、言語学者は自動化された手順でアノテーションを付けられた語彙リストに語彙統計的手法を用いる。そうした手順の一例は子音分類というコンセプトに依拠するもので、ソヴィエト連邦系イスラエル人言語学者のアハロン・ドルゴポルスキーによって1964年に導入された。

P. p b b φ β f v	K. k g x γ q ɠ χ w j	Y. j ç (語根の最初の音の場合)	M. m ŋ
T. t d d θ ð t d	R. r r ɽ ɺ l ɬ k ʎ ʎ t	W. w m (語根の最初の音の場合)	N. n ŋ n ŋ
S. s z ʃ ʒ ʒ z ɛ z c ʃ			Q. t̪ d̪
H. h ʃ h ʃ ʔ h h ʔ, 母音、および j ç w m (語根の最初の音でない場合)			

ドルゴポルスキーの子音分類

以下に世界のさまざまな語族の語彙リストの一部にアノテーションを付けられたものがある。アノテーションは下付き番号とともに与えられている。これらのリストに基づき、*StarlingNJ* アルゴリズムと呼ばれる手法の2つの簡略化したバージョンを用いて語族樹形図が作成され、各単語には安定指数が割り当てられている。上にある樹形図と安定指数は手動でアノテーションをつけられた語彙リストに基づき、下にあるものは自動でアノテーションをつけられた語彙リストに基づく。語彙リストごとに、アルゴリズムAおよびアルゴリズムBの2つのバージョンのアルゴリズムを用いて、それぞれ2つの樹形図が作られている。場合によっては1つの語彙リストに複数の樹形図が対応しうることに注意すること；その場合、ただ1つの樹形図だけがランダムに選ばれている。各樹形図の各ノードには語彙統計的距離が割り当てられている。その距離が大きければ大きいほど、言語間の関係が近くなる。したがって、より正確な用語は「語彙統計的距離」よりも「逆語彙統計的距離」である。簡単のため、この問題では「語彙統計的距離」の用語を用いる。

安定指数および語彙統計的距離は小数点第2位に丸められている。小数点第3位が5未満の場合は切り捨て、それ以外の場合は切り上げ。例えば、2.836は2.84に、0.705は0.71に、0.703は0.70に丸められている。丸めは人間の読者に見える値のみに適用される。すなわち、アルゴリズムを実行するコンピュータは丸められていない値を「見る」。

いくつかの単語は他の言語から借用されたことが判明あるいは推定されているので注意すること。例えば、カディウエウ語の単語 *joki* 「塩」はグアラニー語 *juki* から借用され、イパイ語（メサグランデ方言）の単語 *?a:nj* 「年」はスペイン語 *'ano* から借用されている。

場合によっては、ある単一の意味について、複数の類義語がコマで区切られてその語彙リストの中で与えられることがある。一例はベホス語における「脚」である。

以下のデータにおいて、接頭辞はすべて記号「=」で区切られ、また接尾辞はすべて記号「-」で区切られている。一部の単語は接頭辞なしで用いることができない。これらの単語の前には記号「=」がついている。

データは国際音声字母を用いて転写されている。' = 第1強勢、₁' = 第2強勢（第1強勢よりも弱い）、◌ː = 長い音、◌˚ = とても短い音、XY = XとYは1つの音として発音される、◌ˆ = 高平調、◌˘ = 低平調、◌ˑ = 下降調、◌˚◌ = 前声門化音（のどの空気の通りを短い間止めてから発音

される)、◌' = 放出音 (のどの空気の通りを短い間止めながら発音される)、◌◌ = 無声音、◌◌ = 鼻音化音 (鼻を通して発音される)、◌◌ = きしみ声 (低くパチパチとした音)、◌◌ は子音の前に空気の一部が鼻の中を通ることを示す、◌^h = 有気音 (空気の揺れをともなって発音される)、◌^w = 唇音 (唇を丸めて発音される)、◌^j = 口蓋化音 (舌の一部を硬口蓋に近づけながら発音する)。a、æ、ɛ、ɪ、i、ɔ、u、ɯ、ə、ʌ、ɒ、ɔ、y、θ、ð は母音。その他の特殊な文字は子音。

△ 問題で言及されるいかなる言語の知識もこの問題を解くにあたって有利になることはない。

第1部. グアイクル語族 (アルゼンチン、ブラジル、パラグアイ)

	トバ語 (東部方言)	ピラガ語	モコヴィ語 (チャコ方言)	カディウエウ語
雲	l=ʔok ₁	'lo=ʔok ₁	nawexelek ₂	lol:adi ₃
火	nodek ₁	'd=oleʔ ₂	norek ₁	n=ol:edi ₂
魚	njaq ₁	'nijaq ₁	naʎin ₂	nij:ogo-ɖʒegi ₃
頭	=qajk ₁	'qajk ₁	=qaik ₁	=ak:ilo ₂
殺す	=alawat ₁	=a'la:t ₁	=alawat ₁	=el:owadi ₁
月	ʔawoxojk ₁	ʔa'woʃojk ₁	ʃirajxo ₂	ep:enaj ₃
鼻	=mik ₁	'mik ₁	=mik ₁	=m:iq:o ₁
塩	towe ₁	ol'yek ₂	ʔwe ₁	jok:i ₁
石	qaʔ ₁	'qaʔ ₁	qaʔ ₁	wet:iɡa ₂
舌	=aʈʃ-aʂat ₁	=a'ʈʃ-aʂat ₁	=oʔley-aʂan-aʂat ₂	=ok:eli:ɨ ₃

	アルゴリズム A	アルゴリズム B	
手動	<p>↑ 語彙統計的距離</p>		安定指数: 雲 0.50 火 0.50 魚 0.50 頭 0.75 殺す 1.00 月 0.50 鼻 1.00 塩 0.67 石 0.75 舌 0.50
自動			安定指数: 雲 0.50 火 0.50 魚 0.75 頭 0.75 殺す 1.00 月 0.50 鼻 1.00 塩 0.25 石 0.75 舌 0.50

第2部.ヌビア語族 (エジプト、スーダン)

	ドンゴラウィ語	ケヌジ語	ディリン語	カダル語	デブリ語	ビルギド語
殺す	'bɛ:₁	be:₁	hur₂	wur-i₂	wur-i₂	fila:l-e₁
月	u'n-at-t₁	an-at-t₁	nɔn-ti₁	nɔn-tu₁	nɔn-to₁	ma:l₂
水	'ɛss₁	essi₁	ɔti₁	ɔto₁	ɔtu₁	eji₁
与える	'tir₁	tir₁	ti₁	ti₁	ti₁	te:-n₁
良い	'sɛrɛ:₁	sere:₁	ken₂	kɛn₂	kɛŋ₂	azze-n₃
風	'turug₁	turug₁	irf-i₂	irf-o₂	irf-o₂	kurr-i₃
頭髮	'dil-ti₁	si:r₂	tel-ti₁	til-tu₁	til-tu₁	ur=dill-e₁
お腹	'tu:₁	tu:₁	te-te₂	to₁	to₁	tu:₁
眠る	'nɛ:r₁	ne:r₁	jer₁	dwallɛli₂	jer-i₁	ne:r-i₁
太陽	'masil₁	masil₁	ɛj₂	aju₂	ɛŋgal-to₃	?i:zi₂

	アルゴリズム A	アルゴリズム B	
手動			安定指数: 殺す 0.50 月 0.83 水 1.00 与える 1.00 良い 0.50 風 0.50 頭髮 0.83 お腹 0.83 眠る 0.83 太陽 0.50
自動			安定指数: 殺す 0.33 月 0.50 水 0.50 与える 0.67 良い 0.50 風 0.50 頭髮 0.83 お腹 1.00 眠る 0.50 太陽 0.50

- (A) (2点) 子音**ɛ**はフランス語のrに似た音で、舌の後方で発音される。その子音はどのドルゴポルスキー分類に属するか?どのようにしてその結論を導き出したか?
- (B) (2点) 左上のヌビア諸語の樹形図はあくまで2つある可能な樹形図のうちの一つである。もう一つの可能な樹形図を描きなさい。
- (C) (2点) 左下のヌビア諸語の樹形図はあくまで2つある可能な樹形図のうちの一つである。もう一つの可能な樹形図を描きなさい。
- (D) (2点) (右上のヌビア語の樹形図の根の部分に配置されている) 語彙統計的距離0.49は、この問題におけるその他の距離と同じように、小数点第2位に丸められている。正確な距離はどれくらいか?

第3部. マタグアヤ語族 (アルゼンチン、ボリビア、パラグアイ)

	ウィチ語 (ベルメホ 川下流方 言)	ウィチ 語 (リ バダ ビア方 言)	ベホス 語	ウェーン ナイェク 語	イヨワアハ語	マンファイ語	ニヴァク レ語 (下 流方言)	ニヴァク レ語 (上 流方言)	マカ語
火	ʔitoχ ₁	ʔitəχ ₁	ʔitah ₁	ʔi:taχ ₁	ʰwat ₂	ʔẽitʰe ₁	ʔitax ₁	ʔitax ₁	feʔt ₂
魚	ʔwahat ₁	wahat ₁	wahat ₁	ʔwa:hat ₁	siʔʔjus ₋₁	ʃiʔʔjus ₋₁	saxetʃ ₋₁	saxetʃ ₋₁	sehets ₋₁
脚	=patʃ _{u1}	=qəɓ ₂	=patʃ _{o1} , =kala ₂	=pa:kʔoʔ ₁	=ʰsat ₃	=kaʔʔlaʔ ₂	=φoʔ ₄	=φoʔ ₄	=fʔiʔ ₅
水	ʔinot ₁	ʔinot ₁	wah ₂	ʔina:t ₁	ʔiʔnʔat ₁	ʔaʔnat ₁	jinaʔt ₁	jinaʔt ₁	iweliʔ ₃
与える	=ʔweŋ _{-u1}	=weŋ _{-u1}	=ʔweŋ _{-o1}	=ʔweŋ _{-oʔ1}	=ʰweχn-aʔm ₂	=ʰhajʔ ₃ , =ʰweŋ ₂	=xut ₄	=xut-ej ₄	tis-ix ₅
良い	ʔis ₁	ʔis ₁	ʔis ₁	ʔis ₁	ʔes ₁	ʔẽis ₁	ʔis ₁	ʔis ₁	t=ejkʔun-ej ₂
風	ʔinwok ^w ₁	ʔinwək ₁	ʔihwok ^w ₁	=ja:ʔ ₂ , =x ^w ox ^w ₃	ʰhlahwuʔ ₄	ʰhlahwuʔ ₄	ʔaβiʔm ₅	ʔaβiʔm ₅	tʔunikʔi ₆
木	haʔlo ₁	halo ₁	haʔla ₁	haʔlaʔ ₁	ʔaʔʔlaʔ ₁	ʔaʔʔla-k ₁	ʔaʔkxi-juk ₂	jiʔklaʔ ₁	naxka-k ₃
頭髪	=ʔwule-j ₁	=wule-j ₁	=ʔwole-j ₁	=ʔwo:le-ç ₁ , hi:lenax ₂	=ʔwole ₁	=ʔwole-j ₁	=fateʔtʃ ₃	=jeʔs ₄	=ʔewkux-its ₅
殺す	=lon ₁	=lən ₁	=lan ₁	=la:ŋ ₁	=ʔlaʔan ₁	=ʔlan ₁	=klan ₁	=klan ₁	=lan ₁

	アルゴリズムA	アルゴリズムB	
手動	<p>ウィチ語 (ベルメホ川下流方言) 0.90 ウィチ語 (リバダビア方言) ベホス語 0.80 ウェーンナイエク語 0.44 イヨワアハ語 0.78 マンファイ語 0.33 ニヴァクレ語 (下流方言) 0.78 ニヴァクレ語 (上流方言) 0.11 マカ語</p>	<p>ウィチ語 (ベルメホ川下流方言) 0.90 ウィチ語 (リバダビア方言) ベホス語 0.83 ウェーンナイエク語 0.61 イヨワアハ語 0.78 マンファイ語 0.46 ニヴァクレ語 (下流方言) 0.78 ニヴァクレ語 (上流方言) 0.13 マカ語</p>	<p>安定指数:</p> <p>火魚 0.78 魚脚 1.00 水 0.33 与える 0.78 良い 0.44 風 0.89 木 0.33 頭髪 0.78 殺す 0.67 1.00</p>
自動	<p>ウィチ語 (ベルメホ川下流方言) 0.90 ウィチ語 (リバダビア方言) ベホス語 0.80 ウェーンナイエク語 0.70 イヨワアハ語 0.80 マンファイ語 0.50 ニヴァクレ語 (下流方言) 0.80 ニヴァクレ語 (上流方言) 0.20 マカ語</p>	<p>ウィチ語 (ベルメホ川下流方言) 0.90 ウィチ語 (リバダビア方言) ベホス語 0.80 ウェーンナイエク語 0.75 イヨワアハ語 0.64 マンファイ語 0.80 ニヴァクレ語 (下流方言) 0.80 ニヴァクレ語 (上流方言) 0.28 マカ語</p>	<p>安定指数:</p> <p>火魚 0.78 魚脚 0.44 水 0.33 与える 0.56 良い 0.67 風 0.89 木 0.22 頭髪 0.67 殺す 1.00</p>

第4部. モンゴル語族 (中華人民共和国、モンゴル、ロシア)

(E) (10点) 以下の語彙リストを注意深く研究しなさい。手動および自動のアノテーションそれぞれに対応する安定指数を計算しなさい。

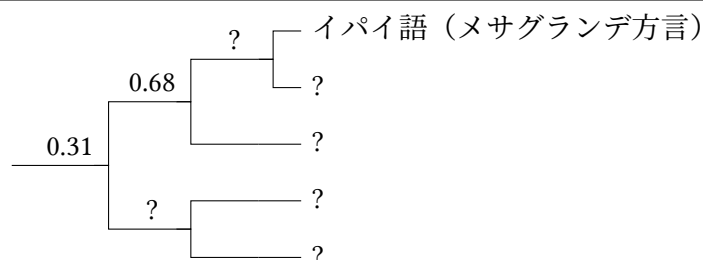
ヒントとして、単語「すべての」については両方の安定指数がすでに計算されている。これらは順不同で0.36、0.40である。

	ダグルル語 (ハイラル方言)	ハムニガト語 (満洲方言)	ブリヤート語 (ホリ方言)	新バルグト語	オイラト語	ホシュト語	カルミク語	ハルハ語	オールドス語	シラ・ユグル語	バオアン語
すべての	hɔ:₁	bolt₂	buxi:₃	bygd₄	tsug₅	lug₅	tsuk₅, xamak₋₁	pux₃, pugt₄, xamäg₋₁	pyyite₄, xamukʰ₋₁	tʃʰuq₅	hanə₋₂
樹皮	hails₁	qalihon₁	χoltōhōn₂	xalʰhu:₁	xolts₂	xalis₁	dursn₃	xəɮtʰɔs₂	turusu₃	χalsən₁	arasun₄
お腹	ke:li₁	getəhōn₂	gedehen₂	gedy:₂	ge:s₂	gets₂	gesn₂	gitis₂, xiwɮʃij₋₁	ketysy₂	ketesən₂	kele₁
鳥	dəgi₋₁	eiwan₁	ʃubu:n₁	ʃuwu:₁	ʃuvu:₁	ʃuwu:₁	ʃowun₁	ʃuwu₁	ʃuβu:₁	ʃu:n₁, peltʃər₂	bendʒer₂
火	gali₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	qal₁	qal₁	χal₁
道	terg-u:l₁	qargūi₂	χargi₂, zam₋₁	zam₋₁	d̄zam₋₁	d̄zam₋₁	xa:-lɤ₃	tsam₋₁	tʃam₋₁	mør₄	mor₄
塩	hata:₁	dawhōn₂	dabhan₂	dawuhu:₂	daws₂	daws₂	dawsn₂	tawsǎ₂	taβusu₂	ta:psən₂	dabsuŋ₂
泳ぐ	unpa-du₁	umba₋₁	tʰamar₋₂	umb₋₁	sele₋₃	umba₋₁	us-t̄ei₋₄, ø:m₋₅	siɮi₋₃	usu-tʃʰi-la₋₄	umpa₋₁	mba₋₁
水	ɔsə₁	ɔxōn₁	uhan₁	u:ha₁	usn₁	us₁	usn₁	ʊsǔ₁	usun₁	qʰusun₁	sə₁
風	kei₁	halkin₂	halxin₂	halxi₂	salʰxin₂	salkʰi₂	salʰkn₂	saɮxi₂	kʰi:₁	kʰi:₁	ki₁

第5部. ユマ語族 (メキシコ、アメリカ合衆国)

(F) (8点) 以下の語彙リストを注意深く研究しなさい。下に、同じ語彙リストに基づいて組み上げられた樹形図がある。いくつかのデータ (言語名および語彙統計的距離) は空欄になっている。空欄を埋めなさい。樹形図が手動生成か自動生成か、およびそれがアルゴリズムAまたはアルゴリズムBのどちらを用いて生成されたか明らかにしなさい。

	モハーヴェ語	ココパ語	ヤヴァパイ語	ティイパイ語 (ハムル方言)	イパイ語 (メサグランデ方言)
短い	wena=wen-a ₁	'xʌ=ʔut ₂	'tʃkr=ot-i ₂	lə=ʔuɲ ₁	mə=put-k ₃
鳥	ʔitʃ=i=jer ₁	'ʃa ₂	'ʔ=tʃ=sa ₂	aʔ='ʃa ₂	ʔa:=ʃa:₂
骨	ɲ=a=s-ak ₁	'ɲ=j=a:k ₁	'tʃ=j=a:k-a ₁	'ak ₁	aq ₁
乾いた	i=ro:-v-k ₁	'ʃ=ʔar ₂	'ru-β-i ₁	's=ʔa:j ₃	sa:j ₃
肉	kʷi:kʷay ₁	ʔi='ma:tʃ ₂	'kʷe:='θo-β-a ₃	'kʷak ₄	kukʷa:j-p ₁
首	maʎaqe ₁	'm=puk ₂	'mlq ₁	i:='puk ₂	i:=puk ₂
見る	i=ju:-k ₁	'wi:₂	'ʔu:₁	'wi:w ₂	ə=wu:w ₂
しっぽ	i:=ʔar ₁	'ʃ=juʎ ₂	'β=hé ₃	ʃə='juʎ ₂	xə=juʎ ₂
二	havik-k ₁	'x=wak ₁	'hʷák-i ₁	xə='wak ₁	xə=wak ₁
年	hu:ðe ₁	'mat-'ka:m ₂	'ʔ=tʃʰur-a ₃	mat-'wam ₂	ʔa:m'₁



(G) (20点) 他にいくつかの樹形図がユマ語族について生成されており、樹形図の根の部分での語彙統計的距離 (すなわち各樹形図のいちばん左側での語彙統計的距離) は以下の通りである:

1. 0.20
2. 0.23
3. 0.24

それぞれの樹形図を描きなさい。各樹形図について、手動生成か自動生成か、およびそれがアルゴリズムAまたはアルゴリズムBのどちらを用いて生成されたか明らかにしなさい。

(H) (3点) 課題 (G) に記載された距離のうち2つが小数点第2位に丸められている: 0.23は0.225から丸められている。もう1つの丸められた距離はどれで、その正確な値は何か?

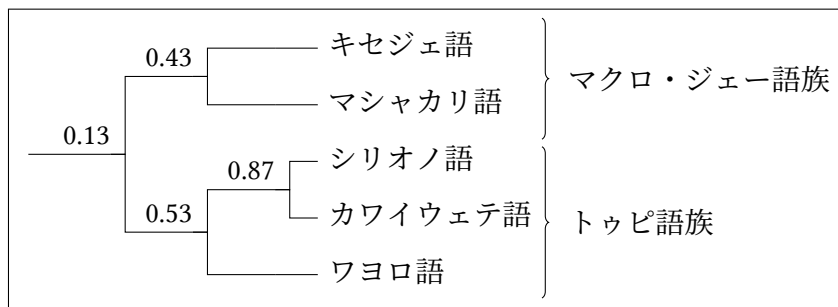
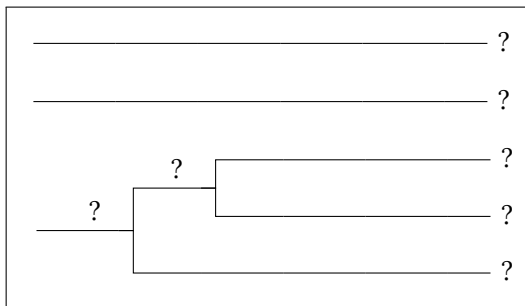
- (I) (4点) 安定指数の計算方法を説明しなさい。
- (J) (5点) 語彙統計的距離の計算方法を説明しなさい。
- (K) (4点) アルゴリズムAとアルゴリズムBの間の違いを説明しなさい。

第6部. マクロ・ジェー語族、トゥピ語族（ブラジル、ボリビア）

(L) (28点) マクロ・ジェー語族とトゥピ語族は南アメリカの2つの主要な語族である。一部の言語学者はふたつの語族が遠縁であると考えている。以下の語彙リストを注意深く研究しなさい。

	A	B	Γ	Δ	E
樹皮	e='e-ke	h ^w i='k ^h Λ	kup='pe	mīβm='təaj	= 'pe
お腹	'e=rje	= 't ^h igi	=ā'ūn	= 'təj	=rε'wek
血	e='ruki	=ka' ⁿ brɔ	=d̄z=a'u	= 'hεβp	=ru'i
燃やす	= 'rai	=rɔ='k ^h Λs̄	=po'k ^w a	mū=...='haβp	=ra'pi
脂肪	e='kira	= 't ^h wəmi	= 'd̄z=ap	= 'tuβp	= 'kap
脚	'e=i	= 'h ^w aji	= 'βi	=po'ta	= 'pi
手	'e=o	=jī' ^k Λa	= 'βo	= 'jīβm	= 'pɔ
重い	e='usi	=wi' ^t hī	=po'ti	=βp'təj	=pɔ'ij
肝臓	'e=ja	= 'nba	=pi'a	=təiβpkī'nāj	=pi'ʔa
新しい	e='jasu	= 'ndiwi	=pa'gop	= 'tiβp	=pia'u
根	e='rao	=ja'ɾe	kup=kujɔ'pe	mīβm=jīβm=təa'tiə	=ra'pɔ
皮膚	'e=i	= 'k ^h Λ	= 'pe	= 'təaj	= 'pit
しっぽ	e='rokoi	= 'nbi	=d̄z=ɔ'k ^w aj	=nā:='kiβp	= 'raj
白い	'e=ʃi	=ja' ^k Λa	=d̄zi'ra	=βp'douɟ	= 'sīj
羽	e='heo	=ja'ɾa	=pe'o	=jī'māuɟ	=pe'pɔ, =ji'wa

下に、同じ語彙リストに基づいて組み上げられた樹形図が2つある。いくつかのデータ（言語名および語彙統計的距離）は空欄になっている。空欄を埋めなさい。各樹形図について、手動生成か自動生成か、およびそれがアルゴリズムAまたはアルゴリズムBのどちらを用いて生成されたか明らかにしなさい。



A	B	Γ	Δ	E
?	?	?	?	?

⚠ この課題では手動アノテーションと安定指数は意図的に省略されている。

- (M) (10点) ドルゴポルスキーの分類に基づき自動化された手順は、誤った結果を出す可能性がある。この例では、自動化された手順は、シリオノ語とトゥピ語族の他言語との間よりも、シリオノ語とある特定のマクロ・ジェー語族の言語（キセジェ語）との間に、より高い類似性を検知する。上のマクロ・ジェー語族とトゥピ語族の語彙リストに適用されたときに正しい分類を出せるよう修正した自動化された手順を提案し、それを簡潔に説明しなさい。

△ この課題は最高得点チームに同点が発生した場合のみ採点される。

著者陣は、アレハンドラ・ヴィダル、マリア・コノシェンコ、イルヤ・グルントフ、ヤムト・スヤに、特定の言語に関する質問について回答してくれたことを感謝する。

—アンドレイ・ニクリン、ミレナ・ヴェネヴァ

編集者: ボリス・イオムディン、スタニスラフ・グレーヴィチ、
イヴァン・デルジャンスキー（技術編集者）、ヒュー・ドブズ、
アンドレイ・ニクリン（編集長）、アレクサンドル・ピペルスキー、
アレクセイ・ペグシェフ、ヤン・ペトル、イーマー・マックナイト、
マリヤ・ルービンシュタイン、エリシア・ワーナー、ミレナ・ヴェネヴァ。

日本語のテキスト: 佐藤和音。

健闘を祈ります！