

فصل سوم : spell correction

۱- الف) سه مورد از کاربردهای استفاده از wildcard query را به همراه مثال بنویسید.

پاسخ :

۱. کاربر در املای یک کلمه از کوئری شک داشته باشد. مثال :

Sydney vs. Sidney, which leads to the wildcard query S*dneý

۲. کاربر بداند یک کلمه با چند املای مختلف نوشته می شود و بخواهد در جستجوی اسناد تمامی حالت های کلمه ای کوئری، برگردانده شود. مثال :

color vs. colour, which leads to the wildcard query col*r

۳. کاربر به دنبال اسنادی باشد که شکل های مختلف یک کلمه را در بر دارند و نمیداند که موتور جستجو ریشه یابی (stemming) را انجام میدهد یا خیر. مثال :

judicial vs. judiciary, leading to the wildcard query judicia*

۴. کاربر از بیان صحیح یک کلمه یا عبارت خارجی مطمئن نباشد. مثال :

the query Universit* Stuttgart

(universitat in German – university in English)

*به هرگونه مثال یا کاربرد صحیح به جز موارد بالا نیز نمره ای کامل تعلق می گیرد.

ب) سه روش زیر که در پردازش wildcard queris استفاده میشود را از نظر سرعت پردازش، قدرت پاسخگویی به کوئری های مختلف و حافظه مصرفی مقایسه کنید.

I. B-tree

II. permuterm index

III. bigram index

پاسخ :

سرعت پردازش : permuterm index > B-tree > bigram index

قدرت پاسخگویی : bigram index > permuterm index > B-tree

حافظه مصرفی : B-tree > permuterm index ~ bigram index

ج) فرض کنید یک لغت نامه شامل کلمات زیر است:

From,romantic,room,chrome, agronomy

نشان دهید این کلمات در دو حالت bigram index , permuterm index چگونه ذخیره می‌شوند. سپس کوئری *rom* و rom* را در هر دو حالت پردازش کنید. تمامی جزئیات را بنویسید.

پاسخ :

الف) روش bigram index :

\$f => from

Fr => from

Ro => agronomy, chrome , from , romantic , room

Om => agronomy, chrome , from , romantic , room

\$r => romantic , room

Ma =>romantic

An => romantic

nt =>romantic

ti => romantic

ic => romantic

oo =>room

\$c => chrome

ch => chrome

hr => chrome

me => chrome

\$a => agronomy

ag => agronomy

gr => agronomy

on => agronomy

no => agronomy

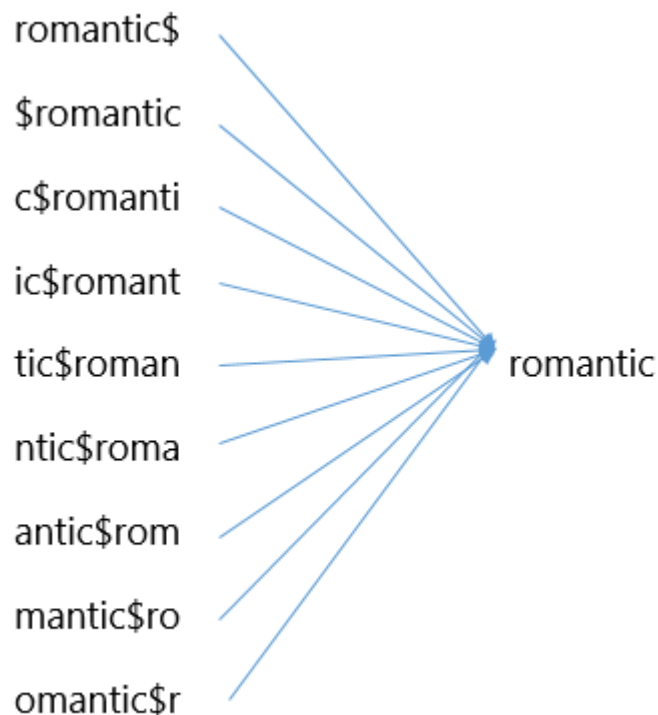
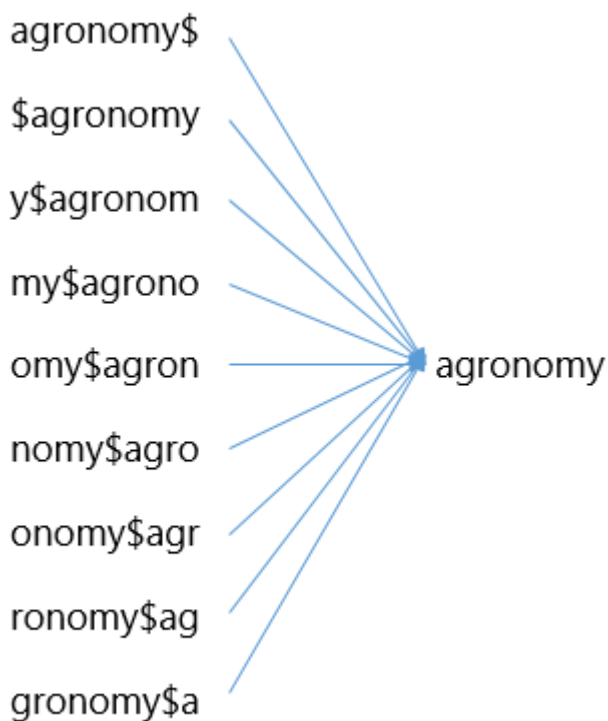
my => agronomy

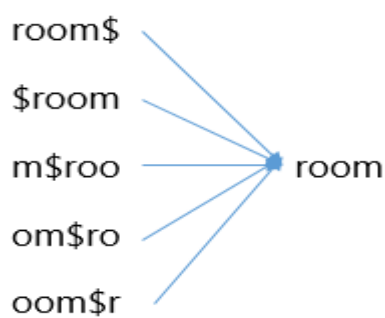
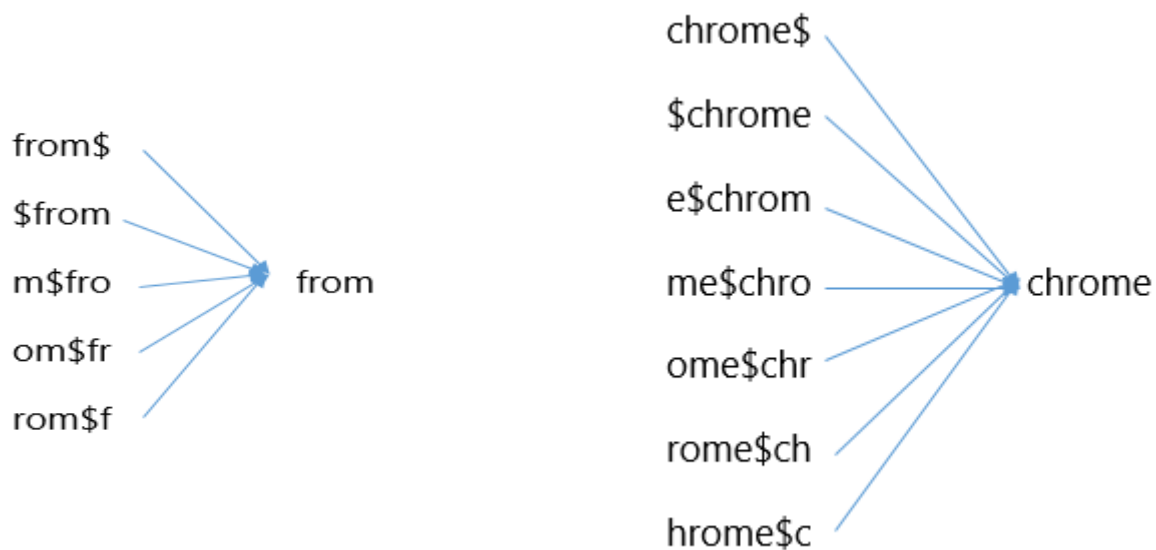
برای کوئری *rom* عبارت ro and om جستجو می شود و حاصل جستجو :

agronomy, chrome , from , romantic , room

برای کوئری *rom* عبارت \$r and ro and om جستجو می شود و حاصل جستجو :

romantic , room





برای کوئری rom^* عبارت rom^* جستجو می‌شود و حاصل جستجو :

chrome , from , romantic

برای کوئری rom^* عبارت rom^* جستجو می‌شود و حاصل جستجو :

romantic

نتیجه :

در کوثری *rom* جوابهای room , agronomy و در کوثری *rom* جواب room صحیح نیستند ولی در خروجی bigram index دیده می شود. بنابراین همانطور که مشاهده می شود در هر دو کوثری عملکرد روش permuterm index بهتر از bigram index است.

۲- در یک موتور جست و جو کلمه "herq" جست و جو شده است ، اما این کلمه در دیکشنری این موتور جست و جو یافت نمی شود . با توجه به پارامترهای داده شده (تعداد رخداد حروف، تعداد رخداد خطاها و احتمال رخداد کلمات مختلف)، بنظر شما کلمه ی وارد شده به اشتباه به جای کدام یک از کلمات زیر وارد شده است ؟

پاسخ:

herq : { "here", "hero", "her" }

تعداد تکرار حرف e : $C(e) = 593086170$

تعداد تکرار حرف r : $C(r) = 32534251$

تعداد تکرار حرف o : $C(o) = 246429812$

$$C(q|e) = 1$$

$$C(q|o) = 1$$

$$C(rq|r) = 1$$

$$P(\text{"here"}) = 0.000624$$

$$P(\text{"hero"}) = 0.00001204$$

$$P(\text{"her"}) = 0.00038$$

Candidate word	correct letter	Error letter	X W	P(X W)	P(W)	$\frac{P(X W)}{P(W) \cdot 10^{10}}$
Here	e	q	q e	$\frac{1}{593086170}$ $=0.00000000016$	0.000624	0.009
Her	-	q	rq r	$\frac{1}{32534251}$ $=0.0000000308$	0.00038	0.1
Hero	o	q	q o	$\frac{1}{246429812}$ $=0.00000000405$	0.00001204	0.0004

با توجه به احتمالات بدست آمده کلمه ی “her” انتخاب می شود.

فصل پنجم : compression

۳- نشان دهید که اندازه لغت نامه در قانون Zipf محدود و در قانون Heaps نامحدود است.

پاسخ :

در قانون Heaps اندازه لغت نامه همراه با افزایش تعداد توکن ها افزایش می یابد، بنابراین می تواند تا بی نهایت پیش برود.

اما در قانون Zipf اندازه لغت نامه باید در ابتدا یک عدد ثابت در نظر گرفته شود، پس محدود است.

۴- الف) در یک سیستم بازیابی اطلاعات یک دیکشنری با ۴۰۰۰۰۰ کلمه داریم، میانگین طول هر کلمه هشت کاراکتر است. اگر بخواهیم این دیکشنری را با روش بلوک بندی و $k=8$ ذخیره کنیم، در حالی که هر اشاره گر به postings list، ۴ بایت و تعداد تکرار هر کلمه (term frequency) در ۴ بایت ذخیره شود و برای اشاره به ابتدای هر بلوک از یک اشاره گر ۳ بایتی استفاده شود، با این شرایط حافظه مورد نیاز را برای ذخیره این دیکشنری بدست آورید .

در این حالت، نسبت به حالت عادی و بدون هیچ گونه فشرده سازی (تخصیص ۲۰ بایت برای ذخیره سازی هر کلمه) چند درصد کاهش حجم حافظه خواهیم داشت ؟

ب) اگر یک دیکشنری شامل مجموعه کلمه زیر داشته باشیم، میانگین تعداد مقایسه هایی را که باید انجام شود تا یک کلمه را پیدا کنیم در دو حالت زیر محاسبه کنید:

۱. بدون فشرده سازی دیکشنری

۲. با فشرده سازی به روش بلوک بندی با $k=4$

{آباد، انقلاب، پیشرفت، تورم، دانشگاه، درخت، رمضان، سلامت، سیاست، فقر، کرونا، وطن}

پاسخ:

(الف)

۴۰۰۰۰۰ کلمه با $k = 8$ ، به ۵۰۰۰۰۰ بلوک تقسیم می شود.

۱. اگر هر کلمه ۸ کاراکتر باشد و یک کاراکتر برای نگه داری طول آن، هر کلمه ۹ بایت حافظه نیاز دارد: 9×400000

۲. هر بلوک به یک اشاره گر ۳ بیتی نیاز دارد: 500000×3

۳. ۴ بایت برای تعداد تکرار کلمات و ۴ بایت برای اشاره گر به postings list: 400000×8

۴. $1 + 2 + 3 = 6950000 \text{ Byte} = 6.95 \text{ MB}$

بدون فشرده سازی:

۱. ۴ بایت برای تعداد تکرار کلمات و ۴ بایت برای اشاره گر به postings list: 400000×8

۲. ۲۰ بایت برای هر کلمه: 400000×20

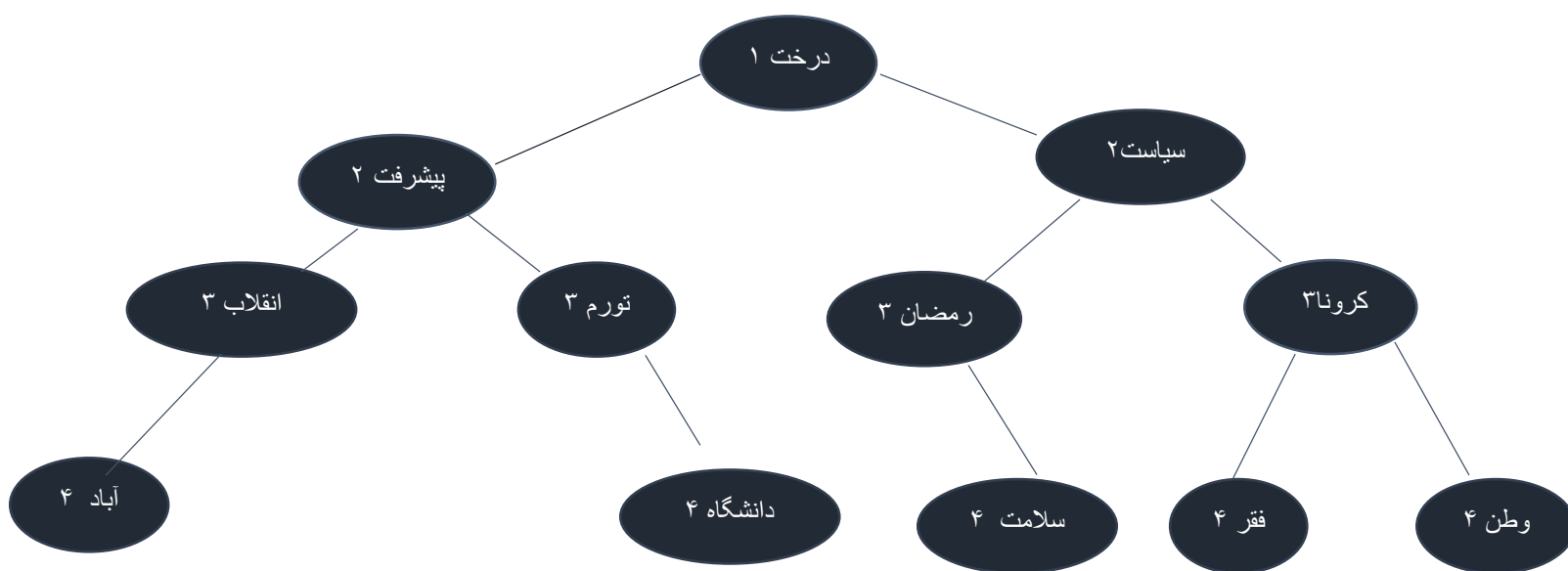
۳. $1 + 2 = 11600000 \text{ Byte} = 11.6 \text{ MB}$

**حافظه مورد نیاز ۴۰ درصد کاهش یافته است.

ب)

{آباد، انقلاب، پیشرفت، تورم، دانشگاه، درخت، رمضان، سلامت، سیاست، فقر، کرونا، وطن}

در حالتی که از روش های فشرده سازی استفاده نمی کنیم، جست و جو بر اساس درخت دودویی زیر انجام می شود.



** تعداد جست و جوی مورد نیاز برای هر کلمه در مقال آن نوشته شده است.

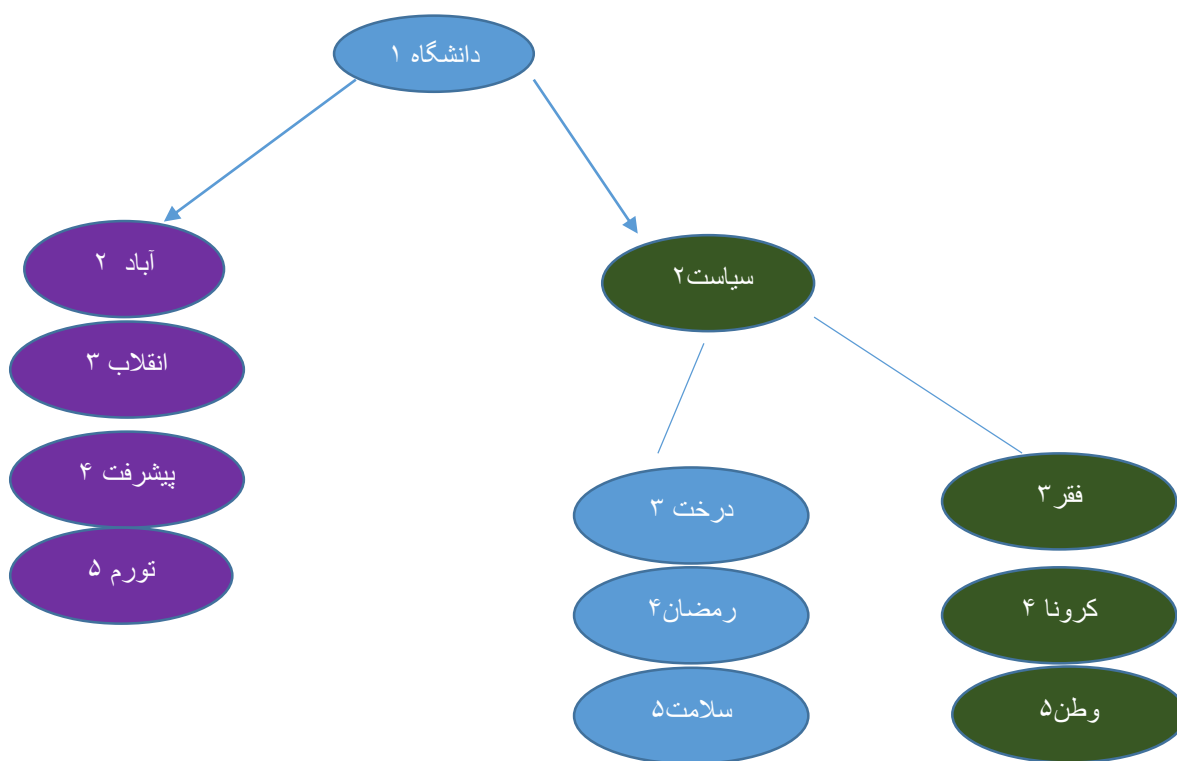
$$(1+2*2+4*3+4*4)/12 = 3.08$$

در حالتی که از روش بلوک بندی استفاده می کنیم :

داده ها به سه بلوک چهارتایی تقسیم می شوند. ابتدا یک جستجوی دودویی بین سه ریشه درخت انجام می شود تا بلوک مورد نظر پیدا شود و وقتی که بلوک پیدا شد داده های بلوک به صورت ترتیبی پیمایش می شوند.

کلمه ورودی با کلمه دانشگاه مقایسه می شود. اگر برابر بود که خاتمه می یابد. اگر کوچکتر بود که لیست سمت چپ به ترتیب بررسی می شود.

اگر بزرگتر بود می بایست یک مقایسه با کلمه سیاست انجام شود. اگر برابر بود که خاتمه می یابد. اگر بزرگتر بود که بلوک سمت راست به ترتیب بررسی می شود و اگر کوچکتر بود که بلوک وسط به ترتیب بررسی می شود.



$$(1*1+2*2+3*3+4*3+5*3)/12 = 3.41$$

۵- یک لغت نامه ی شاخص معکوس شامل کلمات زیر است :

elite ,elope ,ellipse ,eloquent ,eligible ,elongate

الف) فرض کنید دیکشنری در این شاخص معکوس ابتدا کلمات را مرتب و به صورت یک رشته ی طولانی ذخیره می کند. اگر از روش فشرده سازی front coding استفاده شود، رشته ی ذخیره شده به ازای بلوک با سایزهای ۳ و ۶ را بنویسید.

پاسخ :

بلوک سه تایی : $8el*igible3\Diamond ite5\Diamond lipse8elo*ngate2\Diamond pe5\Diamond quent$

بلوک ۶ تایی : $8el*igible3\Diamond ite5\Diamond lipse6\Diamond ongate3\Diamond ope6\Diamond oquent$

ب) فرض کنید :

- در دیکشنری کلمه، تعداد تکرار کلمه و یک اشاره گر به postings list ذخیره می شود.

- تعداد تکرار کلمه و اشاره گر به postings list هر کدام ۴ بایت فضا لازم دارند.

- هر کاراکتر در یک بایت ذخیره می شود.

- برای ذخیره سازی به روش بلوک بندی، هر اشاره گر به کلمه، ۳ بایت فضا اشغال می کند.

- هر کلمه به صورت میانگین ۸ کاراکتر دارد.

- ذخیره سازی با بلوک های ۳ تایی و استفاده از فشرده سازی front coding، باعث کاهش ۲۰٪ از طول رشته ی ذخیره شده می شود.

با توجه به فرضیات ذکر شده سایز بزرگترین لغت نامه ای که می توانیم ذخیره کنیم چقدر است؟ حافظه ی لازم برای ذخیره کردن بزرگترین دیکشنری ممکن چقدر است؟

پاسخ :

پس طول رشته ی کد شده حداکثر می تواند 2^{24} بیت باشد. $\rightarrow 3 \text{ byte} = 24 \text{ bit}$

کلمه $V = 2500000$ $\rightarrow 20 \text{ MB} = V * 8 \text{ B} \rightarrow 20 \text{ M} = 1.25 * 2^{24}$: طول رشته ی واقعی

اشاره گر به کلمه + (تعداد تکرار کلمه + اشاره گر به لیست پستینگ) * سایز بلوک : حافظه ی مصرفی هر بلوک

$$3 \cdot (4 + 4) + 3 = 27 \text{ byte}$$

$2500000 / 3 = 833334$: تعداد بلوک‌های لازم برای ۲۵۰۰۰۰۰ کلمه

16M حافظه مصرفی برای ذخیره سازی رشته‌ی کد شده است. $833334 \cdot 27 + 16M = 38.5 \text{ MB}$: حافظه مصرفی کل
)

۶-الف) کلمه "ایران" در دیکشنری یک سیستم بازیابی اطلاعات در سندهایی با شماره‌های زیر آمده است، این اطلاعات را به روش کدگذاری گاما فشرده کنید.

۲۰, ۵۲, ۷۸, ۱۵۶

ب) اگر postings list را به روش کد گاما ذخیره کنیم، میزان حافظه مورد نیاز برای هر عضو postings list چه قدر است؟

پاسخ:

الف)

docIds : ۲۰, ۵۲, ۷۸, ۱۵۶ \implies gaps : ۲۰, ۳۲, ۲۶, ۷۸

\implies 10100,100000,11010,1001110

number	length	Offset	Gamma code
20	11110	0100	11110,0100
32	111110	00000	111110,00000
26	11110	1010	11110,1010
78	1111110	001110	1111110,001110

(ب)

اگر عدد دهدهی G را به مبنای دو تبدیل کنیم $\lfloor \log_2^G \rfloor + 1$ رقم خواهد داشت.

در قسمت length، به طول معادل دودویی G منهای رقم اول آن، یک و به همراه یک صفر در انتهای آن نگه داری میشود، پس طول قسمت length برابر با $\lfloor \log_2^G \rfloor + 1$ است.

در قسمت offset، عدد دودویی معادل G با حذف پرارزش ترین رقم آن نگه داری می شود، پس طول آن $\lfloor \log_2^G \rfloor$ خواهد بود.

در نتیجه برای عدد G ، $2 * \lfloor \log_2^G \rfloor + 1$ بیت حافظه نیاز است.

۷- در ادامه قسمتی از یک دیکشنری ارائه شده است که postings list آن به روش VB ذخیره شده است، با توجه به این دیکشنری برای پرسمان "تغذیه سالم" چه نتیجه ای به کاربر برگردانده می شود؟

پاسخ:

ابتدا شماره متون را به دهدهی تبدیل میکنیم، میدانیم در این حالت عدد اول شماره متن اول (docId) و شماره های دیگر هر کدام فاصله از شماره متن قبلی (gap) هستند:

تغذیه	10000000(0)-10001100(12)-0000000110001010(138)- 10100100(36)
روان	10000001(1)-10000111(7)-11101010(106)
سالم	10001010(10)-10000010(2)-0000000110001010(138)- 10101000(40)
علم	11000011(67)-10011000(24)-10100011(35)

تغذیه: 0,12,150,186



روان: 1,8,114

سالم: 10,12,150,190

علم: 67,91,26

بنابراین پاسخ پرسمان "تغذیه سالم" متن هایی با شماره ۱۲ و ۱۵۰ است .