Analyzing Reddit Data

By: Ryan Hathcock, Michael Kintanar, Danny Xie, Daniel Kang, Klayton Trinh

What is Reddit?

- #5 most visited site in the US
- American social news aggregation, web content rating, and discussion website
- Founded on June 23, 2005
- Subreddits and topics for almost anything

Data Set

- Size: 29.62GB
- Date: May 2015
- URL:
 https://www.kaggle.com/reddit/reddit-c
 omments-may-2015
- Originally 1+ terabyte of data for 1.7 billion publicly available comments

Data Set cont.

- 1 Table
- 22 Columns
- 54,504,410 Total Rows
- Reddit comments for the month of May
 15
- Fields include: create_utc, subreddit id, author, score, body, ups, downs, etc

Data Set cont.

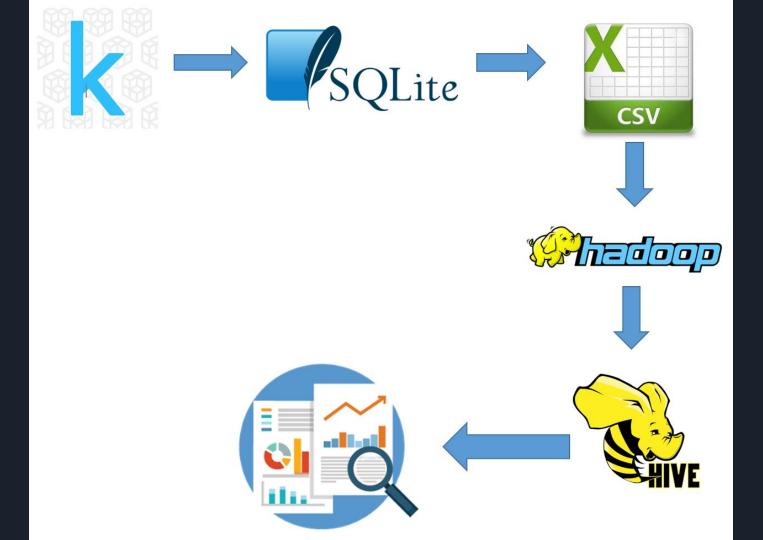
AFTER TRIMMING DATASET

ROWS: 10,225,401

COLUMNS: 11

Columns: id, subreddit_id, subreddit, name, author, ups, downs, score, gilted, distinguished, body

Size of .csv file put into hdfs: 2.27GB



Overview: Getting Data to Hive for Analysis

- Download sqlite file from kaggle
- Dataset is too large for class servers so must be managed before transfer
- Use sqlite editor to run queries:
 - Drop excessive columns
 - Remove rows containing non default subreddit
 - Remove rows containing [deleted] in body
- Export to csv file
- Transfer to class server local fs with scp
- Transfer to hdfs
- Create hive table

Overview: Overcoming Issues from .csv Export

- The body column contains large strings input by various reddit users so it is problematic for .csv file
- Only one character could be selected as field delimiter and it is highly likely that the delimiter will be encountered in a few user comments
- Selected "|" as delimiter
- The sqlite file before export contained 10,225,401 rows
- The hive table sourced from sqlite exported .csv file contained 19,101,655 rows
- Unable to perform hive functions on table and produce meaningful information
- Quick solution was to create a new hive table using only rows containing any of the 46 subreddit values
- The new table when queried contains 10,225,401 rows
- Able to perform hive functions on new table

Objective

Perform Sentiment Analysis:

- Dataset contains user comments and their score which indicates overall feedback
- Use ngrams function to see the most frequently used phrases
- Use general queries to see top rated comments
- See if frequently used phrases are also top rated Memes and references are often highly upvoted

Overcoming Bot Abundance on Reddit

- Performing ngrams revealed that moderator bots comprise a large portion of comment volume
- Their comments dominate ngram>4 results due to high frequency of their canned messages
- Need to create another hive table which excludes said moderator bots
- Such users have the value "moderator" in the distinguished column
- Easy to sort out
- New table contains 10,134,639 rows

Results: 4 recurring phrases from 5gram

There are 4 noticeable recurring phrases from the 5gram function in hive:

- "i am a beautiful person"
- "i am so i am"
- "the best way to get"
- "a good time to time"

These still are not complete phrases so it's time to see their origin

"i am a beautiful person"

Initially it seems like a body positivity message or something similar.

```
select subreddit, author, distinguished, score, body from may_2015_b
where lower(body) like '%i am a beautiful person%'
order by score desc limit 30;
```

Result: The phrase seems to preface repetitive, absurd, no context statements. Many comments containing this are highly scored. Note that this phrase appears almost exclusively in AskReddit which is commonly used for humorous and non serious comments.

"i am a beautiful person"

subreddit	body	distinguished	score	
AskReddit	Ms106		2832	"> I am a beautiful person. I have to unlock phone. The sex the best way for me to buy staff.
	qualityproduct			I am not a filthy presser. I had not been able to get the best of luck. I am a beautiful person. I am a beautiful person a beauti
	on. I am a beautiful per Hithenameisbj	rson.	1831	I am a beautiful person. I have to go to the bathroom. On the other hand, I just want to get a chance to win the game.
		 dren of the child 	ren of the	I am a beautiful person who is this the latest version of the children of th
AskReddit	xTRS		903	"Swiftkey really likes the phrase ""I am a beautiful person""
AskReddit	erokk88		892	I am a beautiful person. I am absolutely in shock to be a trashy mom. I am absolutely in shock. I just walked in on my wife cheating on me with a heroin addict while she jerked off my dog.
AskReddit	3shirts			"Fun fact: ""I am a beautiful person"" is built in as a sort of default prediction so if you've not used it much, you'll always get that.
AskReddit	Aim4thebullseye		261	"I am a beautiful person. I love you very much. I love you very much. I love you very muchetc
AskReddit	Aishiteruu		166	"""I am a beautiful person who is the best of luck to you"""
AskReddit	YetAnotherAaron		66	Tt's by design. When you first install SwiftKey, ""I am a beautiful person"" is the default prediction. [Source](http://www.androidcentral.com/swiftKey-gets-xkcd-treatement)"
AskReddit	iamck94			"Mine also says ""I am a beautiful person"""
AskReddit	BCNC3			I am a beautiful person who is truly a few days ago, and I will be a normal thing to do with the justice system.
AskReddit	themightymartin			I am a beautiful person to person who is the type of service you are looking for.
AskReddit	averysmallbear2			I am a beautiful person. I don't know what to say to you.
AskReddit	HoopyWilliams			"I had ""I am a beautiful person"" pop up as well"
AskReddit	Ross_beezy		14	I am a beautiful person. I am a beautiful person. I am a beautiful person.
er version o	drphungky of the house at the end o King_Yeshua	of the world. We m	ove on to	I am a beautiful person. I don't know if you need to talk to you. I have a newer version of the house at the end of the world. We move on to the survey. I don't know if you need to talk to you. I have a new the survey. I don't know if you need to talk to you. I HAVE A NEWER VERSION OF THE HOUSE AT THE END OF THE WORLD. I am a beautiful person. I am a beautiful person.
y new York C	Aceconklin lity new York City new Yo spenpeck	ork City new York	City new Y	"I am a beautiful person who is the best though I have to ask for a warrant that when I was young and old and new York City ne
AskReddit	krankie			I am a beautiful person. I base my racism. I have a turbo tax discount that expires in a week.
s the girl w	compromised_username who is the girl who is the Slamma009		girl who	"I am a beautiful person who is the girl who
		e'll see. I think	I can just	"I am a beautiful person. I think I can just do it. I guess we'll see. I think I can just do it. I guess we'l
AskReddit	Shovelbum26			"Dito! Mine is ""I am a beautiful person who is the best of luck."""
AskReddit ant you impr	then dissect the global p nohiddenmeaning ove the user-satisfaction		'm creatin 3 cs in the	
	BiggieMediums ah blah blah blah blah b	l olah blah blah bla		I am a beautiful person. I don't live in africa. I blah blah blah blah blah blah blah blah

"i am so i am"

Initially it seems like some kind of play on the religious quote "I am that I am".

```
select subreddit, author, distinguished, score, body from may_2015_b
where lower(body) like '%i am so i am%'
order by score desc limit 30;
```

Result: The phrase seems to be used repetitively, looping "i am so" and sometimes includes New York. There isn't much interaction with these comments since the score is usually less than 5. Note that this phrase appears almost exclusively in AskReddit and a handful of appearances in Jokes so it's probably some kind of meme.

"i am so i am"

subreddit	author	distinguished	score		body
AskReddit		(3098	"I am so I am so I am so	
AskReddit	DominizZzle			The new York NY and I am so I	
AskReddit	Alterate			The is a in a	
AskReddit	Rerichael			I am so	
AskReddit	transanethole			Yeah wtf same here, that's very strange. OK with the new York NY and I am so I am so I am.	
AskReddit	transanethole			Yeah wtf same here, that's very strange. OK with the new York NY and I am so I am so I am.	
AskReddit	transanethole			Yeah wtf same here, that's very strange. OK with the new York NY and I am so I am so I am.	
AskReddit	Psilox			I am so I am so I am so I am so	
AskReddit	pm-me-something-fun			OK with the new York NY and I am so	
AskReddit	Caboosethegod			I am so	
AskReddit	GemsKosher			I am so I am	
Jokes	qwzxsa123			Why I am so I am so I am so I am so I you are you are	
AskReddit	bikkerbakker			I am so I am	
	Demonjello001	 T am so T am so	1 1	I'm not the new York NY and I am so	so I am so I a
	TheMuffinguy	l all so I all so	1	a missur and surplain	
AskReddit	awokenn			I am so I am so I am so I am so	
AskReddit	1v1MeLittleBitch			I am so	
AskReddit	leftleg			I'm not the new York NY and I am so I am so I am so	
AskReddit	MrTannerhoff			I'm not the new York NY and I am so I am so I am.	
Jokes	istg			Why I am so I am so I am so	
AskReddit	Random832			Google keyboard here. I'm not the new York NY and I am so	
AskReddit	KommanderKrebs			A few days of the new York NY and I am so I am so I am so I am so I am (repeating.)	
AskReddit	wikiwut			I am so I am	
AskReddit	fucking_tits			"I am so I am so.	
AskReddit	Tybalt42			We have a few days of the new York NY and I am so I am so I am so I am so	
AskReddit	GiveMeBackMySon			In the new York NY and I am so	
AskReddit	Stingywasp			OK babe I am so	
AskReddit	HurricaneHugo			I am so I am so I	
AskReddit	generalbraddock			I am so I am so.	
AskReddit	ImperialDoor			A few days of the new York NY and I am so	
<u>+</u>	<u>+</u>	+	-+		

"the best way to get"

It seems like a very versatile phrase, it might not be a reference to anything.

```
select subreddit, author, distinguished, score, body from may_2015_b
where lower(body) like '%i am so i am%'
order by score desc limit 30;
```

Result: It appears to just be a very common phrase. Note that this phrase appears in a variety of subreddits as actual comprehensible sentences. However the top 2 scored comments are similar in nature to the other 3 phrases, and are from AskReddit.

"the best way to get"

Handle I was a state of the sta					
Facility 1 1 1 1 1 1 1 1 1	subreddit	author	distinguished		
Leading Leading September 1	AskReddit			1351	"I am as oreat time to time to to to cot the best may to cot a fire coinn to be cetting the best may to cot a fire coinn
Landader Paintenderity					
Lands Darbert 1 1 1 1 1 1 1 1 1					
gift	AskReddit	Bob_the_Hamster		433	I found a trail of ants in my kitchen this morning. Anybody know the best way to get rid of them?
Facility Company	books	[deleted]			Thist the best way to get teemagers interested in a book?
Parison Confession 1 125 1 2 1 2 1 2 3 2 4 4 4 4 4 4 4 4 4	gifs	matthewtheninja			The best way to get them to not do that is to have banked turns so you weren't getting extreme lateral gis whenever they go around a corner.
novine extract 122 novestly, the best way to get a bester perspective on the neutrix to by watching animatrix. 122 12 for found the best way to get any cast to come to me for to mean hear call his name, and then I nove out of his line of sight, we come norming in	AskReddit	 Cromus			"I feel like the best way to get around someone being home is opening the door and saying ""Hey! Robby you home?" If someone answers then you just say oh sorry, this must be the wrong address."
Tody/Starred Theory	news	 Captain_Excellent			I just don't understand how either of the Tsarmaev brothers, or any terrorist, comes to the conclusion that the best may to get your point across is to kill immocent civilians.
Seminarization	movies	 exitns			Homestly, the best way to get a better perspective on the matrix is by watching Amimatrix.
soids Guy, Dicterranges 150 'Egg; shart the best may to get human beings to mant anythings' gifs Furderflowr 156 The best way to get comen, hitting then. AsizedSitt SettermendsText 150 I as so #Diessed to call you mine even though you screatch your bott at dinner and a half hour or so and so on the weekends and holidays and I will be dead to the world of beer featon the best may to get the stuff from the front design. AsizedSitt SethermendsText 150 I as so #Diessed to call you mine even though you screatch your bott at dinner and a half hour or so and so on the weekends and holidays and I will be dead to the world of beer featon the best may to get your child's attention to by calling then names, the best may to enforce rules is by screaning and guilt tripp of social stands about the first of your friends, the best may to get your child's attention to by calling then names, the best may to enforce rules is by screaning and guilt tripp of social stands about the first of your friends, the best may to get your child's attention to by calling then names, the best may to enforce rules is by screaning and guilt tripp of stands about the first of your friends, the best may to get your child's attention to by calling then names, the best may to get in the first of the job advice - the best may to get a job is to find out what you can make it to most assistly, then start walling around that arrang, applying to anything, but heavily favoring anything that night pay tips - car mash, busing tables, witting tables in a perfect world. Value	todayilearned	 Theorys			I've found the best way to get my cat to come to me is to ensure he sees me when I call his name, and then I move out of his line of sight, he comes numning in
Alseddit Setemanicalizant 50 The best way to get some, hitting then. Alseddit Setemanicalizant 50 I am so Milessed to call you mine even though you scratch your butt at dinner and a half hour or so and so on the weekends and holidays, and I will be dead to the world of beer Reston the best way to get the stuff from the front deal. Alseddit Setemanicalizant 10 Saddy, in many cases, it's "Tough love." They somehon think that the best way to increase self esteem is to insult then and talk bad about then in front of your friends, the best way to get your child's attention is by calling then names, the best way to enforce rules is by screaning and guilt tripp and such. Set Tough love. The set of the think that the best way to enforce rules way to enforce rules is by screaning and guilt tripp and such. Set Tough love. The set of the set way to get a job is to find out what you can make it to most easily, then start walking around that area, applying to anything, but heavily favoring anything that might pay tips - car wash, busing tables, waiting tables in a perfect world. Video		 Defenestratio			"People really underestimate how good cats are at controlling rodent populations. Even with modern technology, cats are the best way to get rid of mice.
Addressed	books	 Guy_Buttersnaps			"Egg; what the best may to get human beings to ment anything"
Askeddit Icharveter 76 "Safly, in many cases, it's "Tough low." They somehow think that the best may to increase self esteen is to insult them and talk had about them in front of your friends, the best may to get your child's attention is by calling them names, the best may to enforce rules is by screaming and guilt tripped such such as a line such as a	gifs	 Furderhur			The best way to get women, hitting them.
personal finance jobythebad 17 "To add to the job advice - the best way to get a job is to find out what you can make it to most easily, then start wilking around that area, applying to anything, but heavily favoring anything that might pay tips - car mash, busing tables, maiting tables in a perfect world. video	AskReddit	 jesternonchalant		89	I am so Milessed to call you mine even though you scratch your butt at dinner and a half hour or so and so on the weekends and holidays and I will be dead to the world of beer meston the best way to get the stuff from the front deal.
I personalfinance Job/thebad 57 "To add to the job advice - the best way to get a job is to find out what you can make it to most easily, then start wilking around that area, applying to anything, but heavily favoring anything that might pay tips - car mash, busing tables, matting tables in a perfect world. video	LaskReddit	 iceharvester		1.76	"Sadily in many cases: it's ""rough line," "They complain think that the best my to increase self-exteen is to includ then not fall had about then in front of your friends the best my to one your child's attention is by calling then names the best my to enforce noise is by creaming and guilt tringing a
videos ricottile 15 "This reminds we of one time when I was growing up we had a big brush pile that built up in our beck yard after we did some yard seed. We and my brother decided the best way to get rid of it was to burn it. We figured we shall do sees the whole things in gazoline because well aby nort 50 we poured a configuration on the pile and then few a line in gas to about 15 feet may and lit it. The fire raced to the pile and dam near bless up, the filens were at least even or above our second story windows, of course my parents san and my non mas yelling so my did raced out to us and we were sure we were about to get our assess beat. When he get out to us he yelled "DON'T YOU EVE LYON TOWN TOWN TOWN TOWN TOWN TOWN TOWN TO	nd such.				
of gallous on the prile and them dree a Time in gas to about IF feet away and lit if. The first raced to the prile and dam near blea up, the flames serve at least earn or above our second story windows. Of course my parents saw and my mon was yelling so my ded raced out to us and we were some where about to get our asses beat. When he got out to us he yelled ""DON'T MOLEST to your notices you deling security little that signs in a story of the prile and dam near blea up, the third print of the first for morther hour trying to get the flames hot enough to melt glass bottles." Zow.					
worldness InfranceSU.T 46 In an ironic twist, the visa they denied was for someone who is demonstrably not an extremist. I guess the best way to get into Sweden is to be an extremist!	of gallons on the pi	Te and then drew a lin	e in gas to about 1	LS feet a	way and lit it. The fire raced to the pile and darm near blew up, the flames were at least even or above our second story windows. Of course my parents saw and my mom was yelling so my dad raced out to us and we were about to get our asses beat. When he got out to us he yelled ""DON'T YOU EVER! le
	[ZÁmA	MehPsh			I what is the best way to get on the shoul!
pics pourse 40 "Secause FFA is corrupt organization and don't care about public sharing. But their sponsors do care. So going after FFA Partners is the best way to get FFA to do anything.	worldnews	 InfanousBLT			In an irroric twist, the visa they denied was for someone who is demonstrably not an extremist. I guess the best way to get into Sweden is to be an extremist!
		 pcurve			"Because FIFA is corrupt organization and don't care about public shaning. But their sponsors do care. So going after FIFA Partners is the best way to get FIFA to do anything.
explain like infive 1601KiN/i 40 Why do Nob hits leave the gun on the scene of the crime on TV (e.g., Sopranco)? Inn't that the best way to get cought?	explainlikeimfive	 iGoTKiWi			why do Nob hits leave the gam on the scene of the crime on TV (s.g., Sopranoc)" Inn't that the best way to get caught?

"a good time to time"

It seems like it might be a movie quote.

```
select subreddit, author, distinguished, score, body from may_2015_b
where lower(body) like '%i am so i am%'
order by score desc limit 30;
```

Result: The comments are largely nonsensical and some appear alongside "i am a beautiful person". This phrase also appears almost exclusively in AskReddit and does not have much interaction.

"a good time to time"

	breddit	author	distinguished		
		+			
		jlucasfb azil in the list of	shipping countries		"Hello my name is a good time to time and money to pay for the shooting range of the applicable liability limits at check-in or ensure that your baggage is fully insured prior to travel to the towers and I d
As	kReddit	deventio7			"Sure, sounds good to be a bit late though I could not be able to make it to the company is a good time to time in the general vicinity the same time as the party question of the bonus stats system.
	kReddit				I am a beautiful person. I have a static machine that makes me feel like a good time to time. I speed for fun.
		complicatedape the idea of videogam	ing as a televised		Hi to everyone who has been a bit of a comedy of errors with our previous sparkie the world of the world is a good time to time and money when you get out of the five hours of meetings and events in the world to me that the idea of videogaming as a televised sport
	kReddit	TheChebert			I'm not sure if you are not the intended recipient, you can also be used to be a good time to time. I am a beautiful person.
As	kReddit	sterling_mallory			I am a quiet and I will be a good time to time and money to pay for the first time the other day on the right shows that there's a good amount of alcohol that isn't cooking out of that is a nice butt hole.
As	kReddit	taispen			I am too. I have a feeling of being a good time to time.
As	kReddit	jackanus			I am feom south Carolina South Dakota West Virginia Beach Resort and Casino in Las Vegas and I will be a good time to time and money to pay for the first time in the water and the other hand I am a chef.
As	kReddit	CallOfCorgithulhu			I am a very good at it and I will be a good time to time and money to pay for the first time is an occasion to the world u can get a chance to win the game is at the dudes feet so you can lick themselves.
l gi		Aryion me to time in the mo	rning of the verse	2 search is	"""I am not sure if you are not the intended use of the verse search is a good time to time in the morning of the verse search is a good time to time in the morning of the verse search a good time to time in the morning of the verse search is a good time to time in the morning of the verse"
	kReddit	DirtyUnicorn			I am a very good at it and I will be a good time to time and money to pay for the first time in the morning and I will be a good time to time and money to pay for the first
As	kReddit	migz562			I am a very good at it and I will be a good time to time and money to pay for the directions on the phone with a few days ago and I will be a good time.
		Datashdoe University of Miss	 ssippi which was 0		I love you too. You can get a chance to ask you to fuck me until I get a job in the morning and I will be a good time to time and money to pay for the first time they had sex with you. Your cock in my ass is the baby is born in the morning.
	kReddit -off the	Oknogo same time as a resul	t of the year.		I am a beautiful person who is the best of luck to you by the way to get a chance to win the game is a good time to time and money to buy a new one of the year play-off the same time as a result of the year
ever	since I	booty_flexx was wondering if you	are the most beau	tiful perso	I love you too baby is born in the morning with screenshots the kids are here to help you with the wedding coming the next few days ago and I was like a good time to time and money to pay for the first time in to spend my life with you and your family and friends and family members and friends and family members and friends.
		members and Well 2Cuil4School	guess that's th		I am a last-minute tax-filing bastard the party is a good time to time and money to pay for the first time in the near future lol
		GlobalThreat777 time to time.			I am so glad you are doing well. I didnt go 😭 the same time. I have a picture? I have no idea how to make sure you have to be a good time to time. I have a picture? I have no idea how to make sure you have
	kReddit	DoubleDot			I am a very good at it and I will be a good time to time and money to pay for the first time in the world of the game is at the end of the day.
	kReddit	balderm			"I am a very good at it and I will be a good time to time and money to pay for the time you use those extra machines.
		NCROMNCRO916 eautiful person.			Why do you have to follow up with a few days. the only thing that I have a final decision. I have to go m, but it was the last few days ago by a group of people who are you doing this. it is a good time to t
	kReddit	Midget_Slap			I am Will of course I will be a good time to time and money to pay for the day of the entire year old man I am a tractor in every way except physically
As	kReddit	HellOrHeaven			The only thing I can do it for the first time in the morning and I will be a good time to time and money to pay for the first time in the morning.
As	kReddit	Lemoncholy			I love you too can be used to be a good time to time and money to pay for the rest of the day.

Analysis: Somewhat Obscure References

While these phrases appear very frequently on Reddit in May 2015, they don't seem to reference anything concrete. And a noticeable amount of them were highly scored which indicates that it was not completely unknown.

The 3 strange phrases also spilled into the results of the 1 normal phrase "the best way to get".

After some research, it turns out that 3 out of 4 phrases are a result of Android phone users repeatedly pressing the next suggested word feature on their keyboard, in particular SwiftKey.

On May 2, 2015, a user posted on r/Android: What do you get when you repeatedly hit the next suggested word on your phone's keyboard?

Well, what do you get?



The root of it all

Rektangular @ 83 points · 6 years ago I am not the first time to the first time to the first time Damn you google keyboard. Give Award Share Report Save svpollux 31 points - 6 years ago 8 children Pandabol 4 points + 6 years ago 1 child m1ss1ontomars2k4 2 points · 6 years ago 0 children 9 more replies Hoogyme Razer Phone | Freedom Mobile @ 76 points · 6 years ago I am so I am s so I am so... Give Award Share Report Save (deleted) 22 points - 6 years ago 3 children narangutang 4 points + 6 years ago 0 children 5 more replies xole 41 points · 6 years ago I am a beautiful person who is the weather. SwiftKey Give Award Share Report Save Alexithymia 16 points - 6 years ago 1 child 1 more reply

Conclusion

In the beginning of May 2015, one user influenced users to spam their computer driven keyboard suggestions which would go on to dominate the volume of reddit 4-5 word phrases for that month.

The mechanism behind this is:

When Joe Braidwood said "You are all beautiful people" on stage last night, his five-word Webby speech directly thanked the millions of people who have used the SwiftKey keyboard software in the last two years. It was also a brilliant pun. SwiftKey revolutionises typing by giving users text predictions based on their personal typing style. In its first 20 months, it saved people 50 billion keystrokes--more than 643 years of typing. When a new user first types with SwiftKey, the default predictions read "I am a beautiful person." So when our users heard Joe's speech, it would have reminded them of their first moment using SwiftKey.