



Let's Play a Game  
Never Have You Ever.....

# Never Have You Ever Had an 8-year Relationship



2015  
grade 8



2018  
grade 11



2021  
university

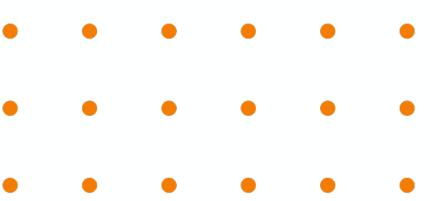
# **Why Did I Broke Up with My Ex after 8 Years**

--- a Pattern and Sentiment Analysis  
Based on the Chat History



Presented By:  
Cheng Ling Jun  
3035772652

# 1. Problem



## Why Using Data Science?

- Emotion Leads to One-sided Opinions, especially in a relationship
- Facts that could be misunderstood due to one's emotions and perspective
- Data science research is needed to find out the quantitative and objective “truth”.

## Expected Outcome?

- Pattern
  - Evolution of the chat frequency
  - Data insight on an monthly/weekly/daily basis
- Sentiment
  - the effect brought by positive and negative messages
  - Comparison of average sentiment value over the control group



# 1. Problem



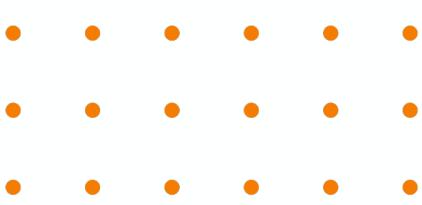
## Why Important?

- If you ever had an 8-year-old relationship.....



## 2. Dataset

- Data Source:
  - WeChat History, 2021 May - 2024 April



localId	TalkerId	Type	SubType	IsSender	CreateTime	Status	StrContent	StrTime	Remark	NickName	SenderId	Sender
22855	28	1	0	1	1620014346		你在干啥呀	#####	Clara Cheng	Clara Cheng	wxid_	[REDACTED]
22856	28	1	0	0	1620015553		刚刚吃完饭	#####	[REDACTED]	[REDACTED]	[REDACTED]	[REDACTED]
22857	28	1	0	0	1620015557		准备洗碗	#####	[REDACTED]	[REDACTED]	[REDACTED]	[REDACTED]

- Tools



Memotrace  
to conduct data scrapping  
from WeChat



R Studio  
for pattern analysis



Weiciyun  
for sentiment analysis  
(understanding Chinese Better)

## 2. Dataset



- Data Cleaning:
  - String Only
    - Only string content is saved.
    - Other types of data is not included, such as picture, video, emoji, etc.,
  - Only after 2021 May
    - *CreateTime* is a number defined by Memotrace, which represents the time



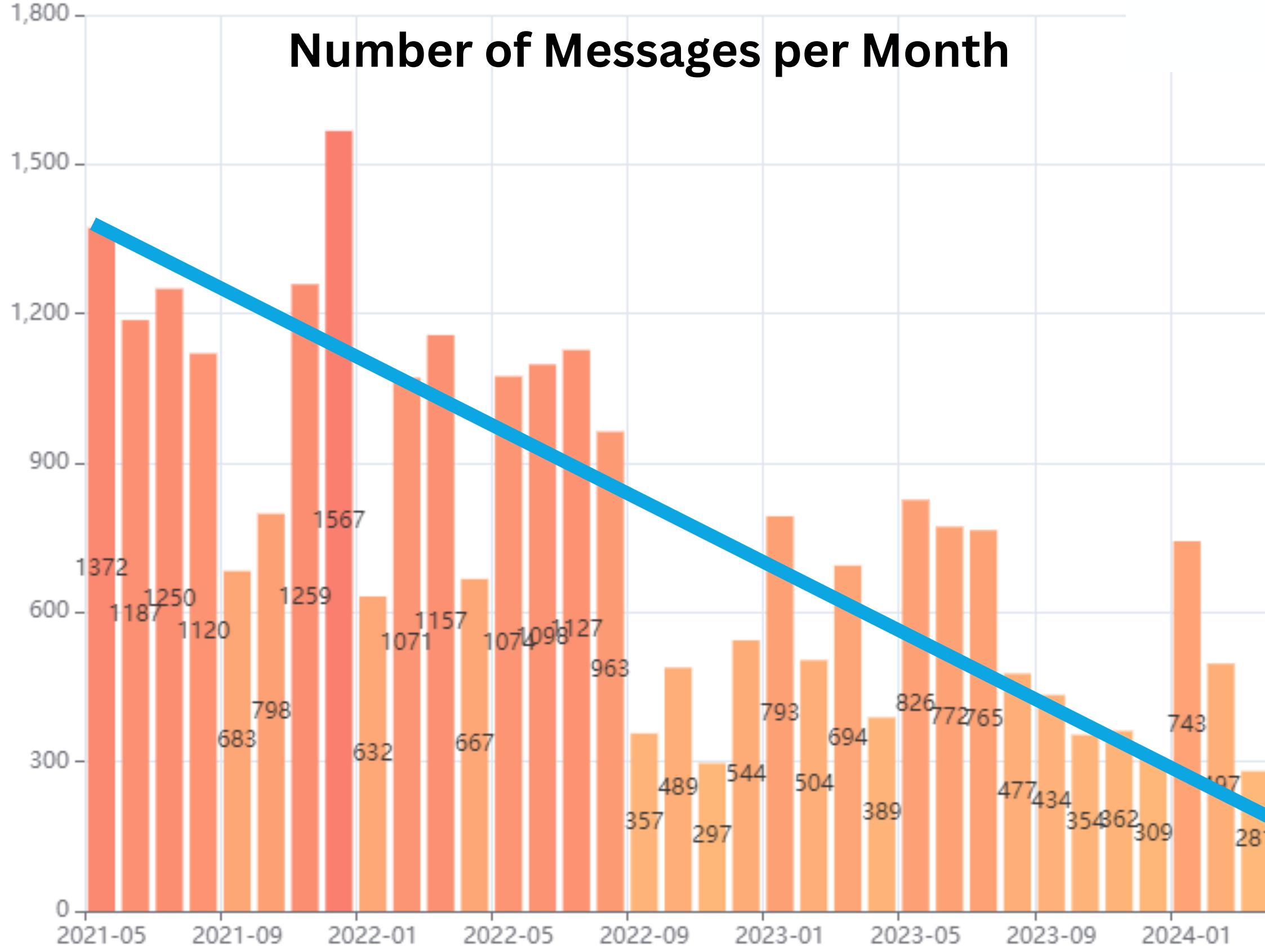
Filter

Goal	Condition
String Only	<i>Type</i> == 1
After 2021 May	<i>CreateTime</i> > 1620014346

localId	TalkerId	Type	SubType	IsSender	CreateTime	Status	StrContent	StrTime	Remark	NickName	Sender	
22855	28	1	0	1	1620014346		你在干啥呀	#####	Clara Cheng	Clara Cheng	wxid_	
22856	28	1	0	0	1620015553		刚刚吃完饭	#####				
22857	28	1	0	0	1620015557		准备洗碗	#####				
22858	28	47	0	0	1620015558		<msg><emoji>	#####				

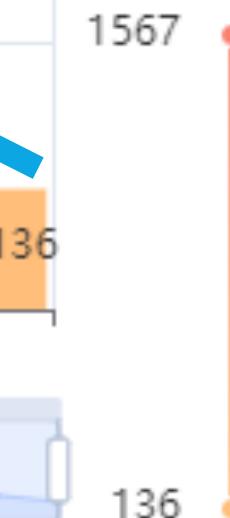
Filtered out

### 3. Pattern Analysis - Month

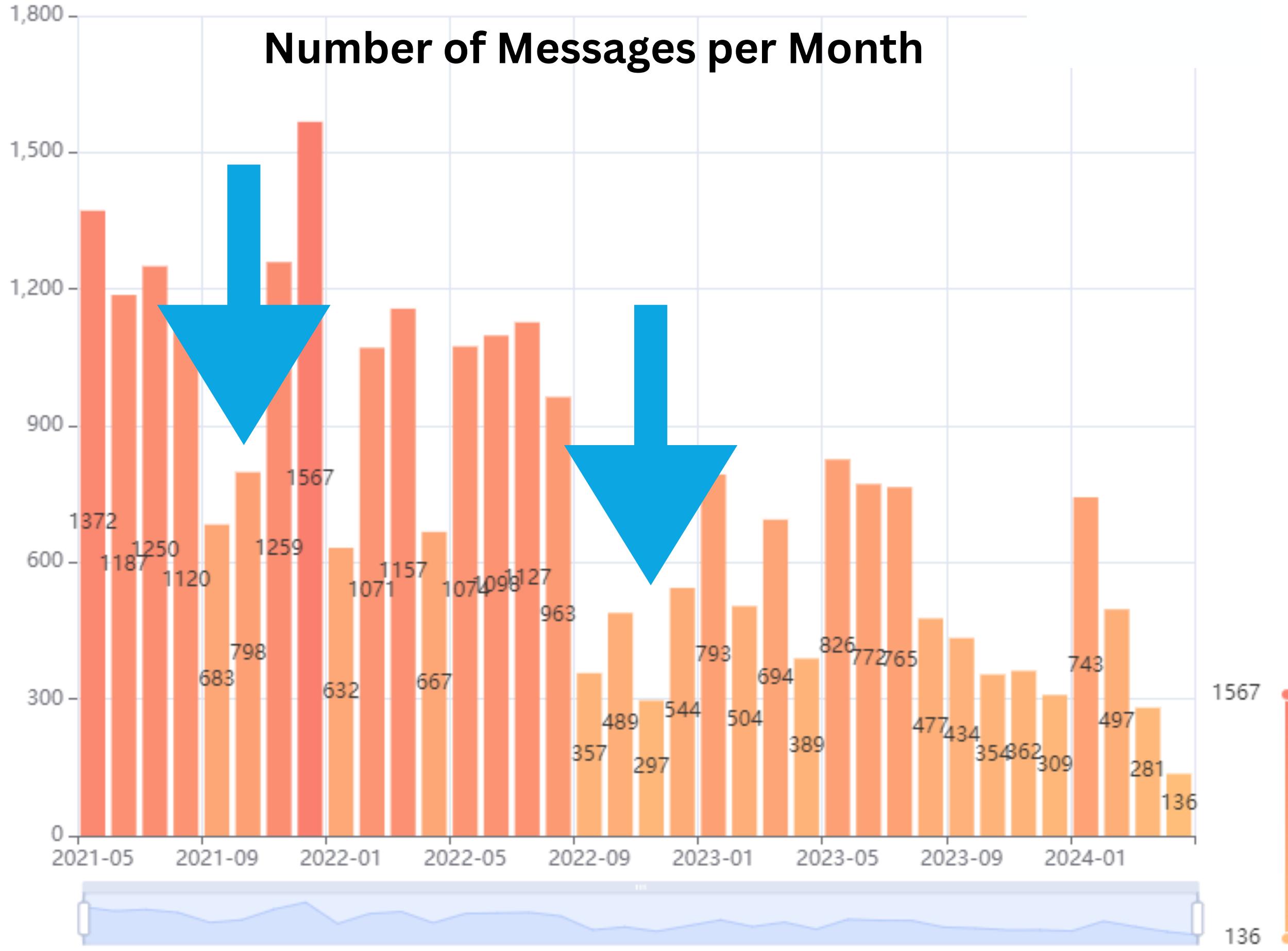


Observations:

- The number of messages decreased steadily by 90% in the past 3 years

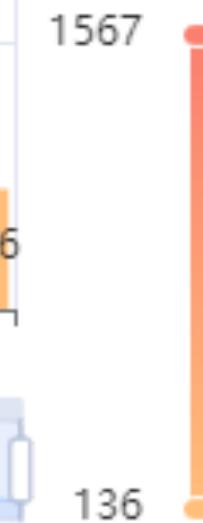


### 3. Pattern Analysis - Month



Observations:

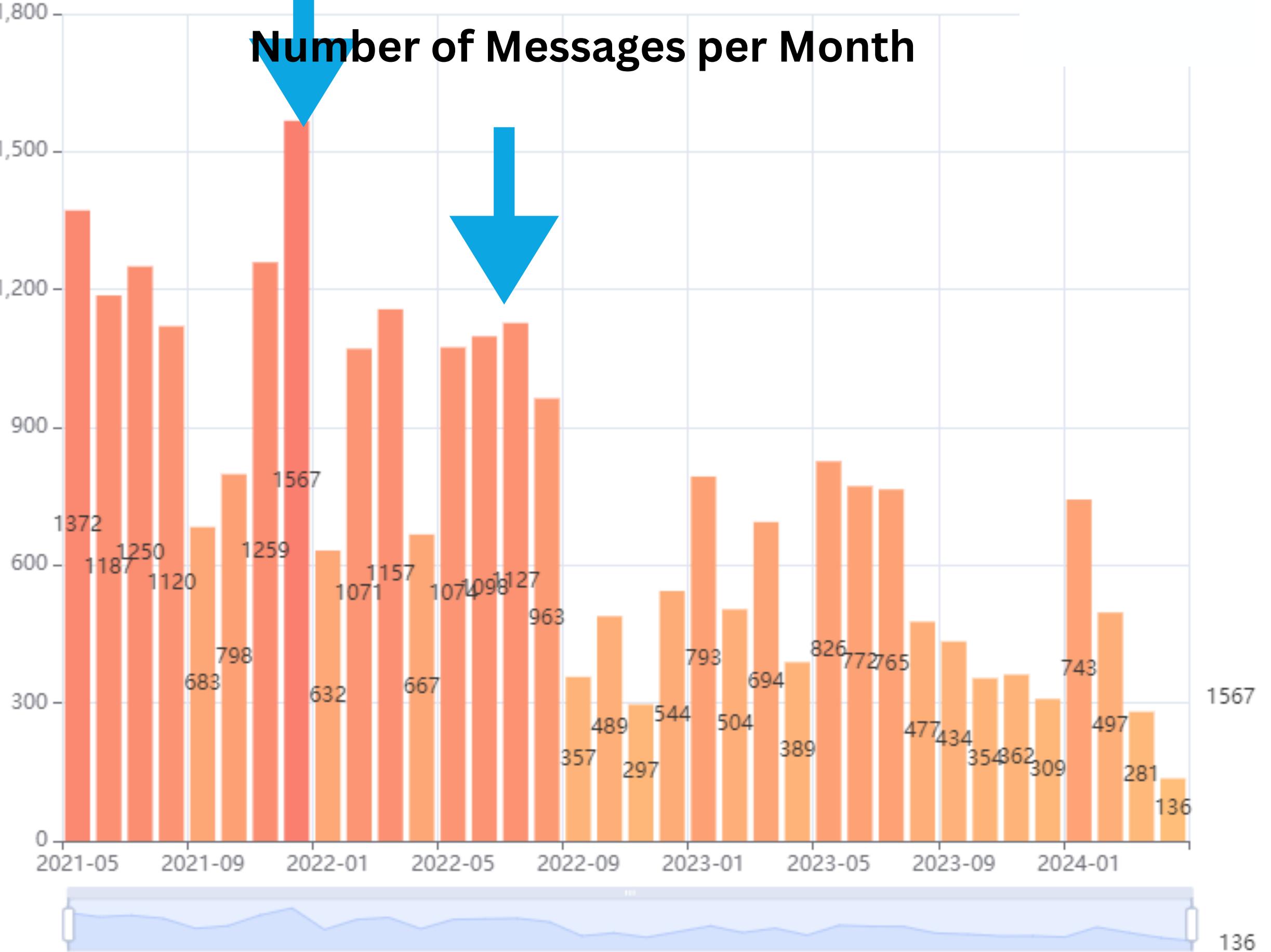
- The number of messages decreased steadily by 90% in the past 3 years
- less chat at the beginning of a semester instead of “blood final”



### 3. Pattern Analysis - Month



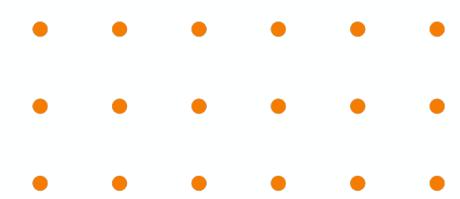
Number of Messages per Month



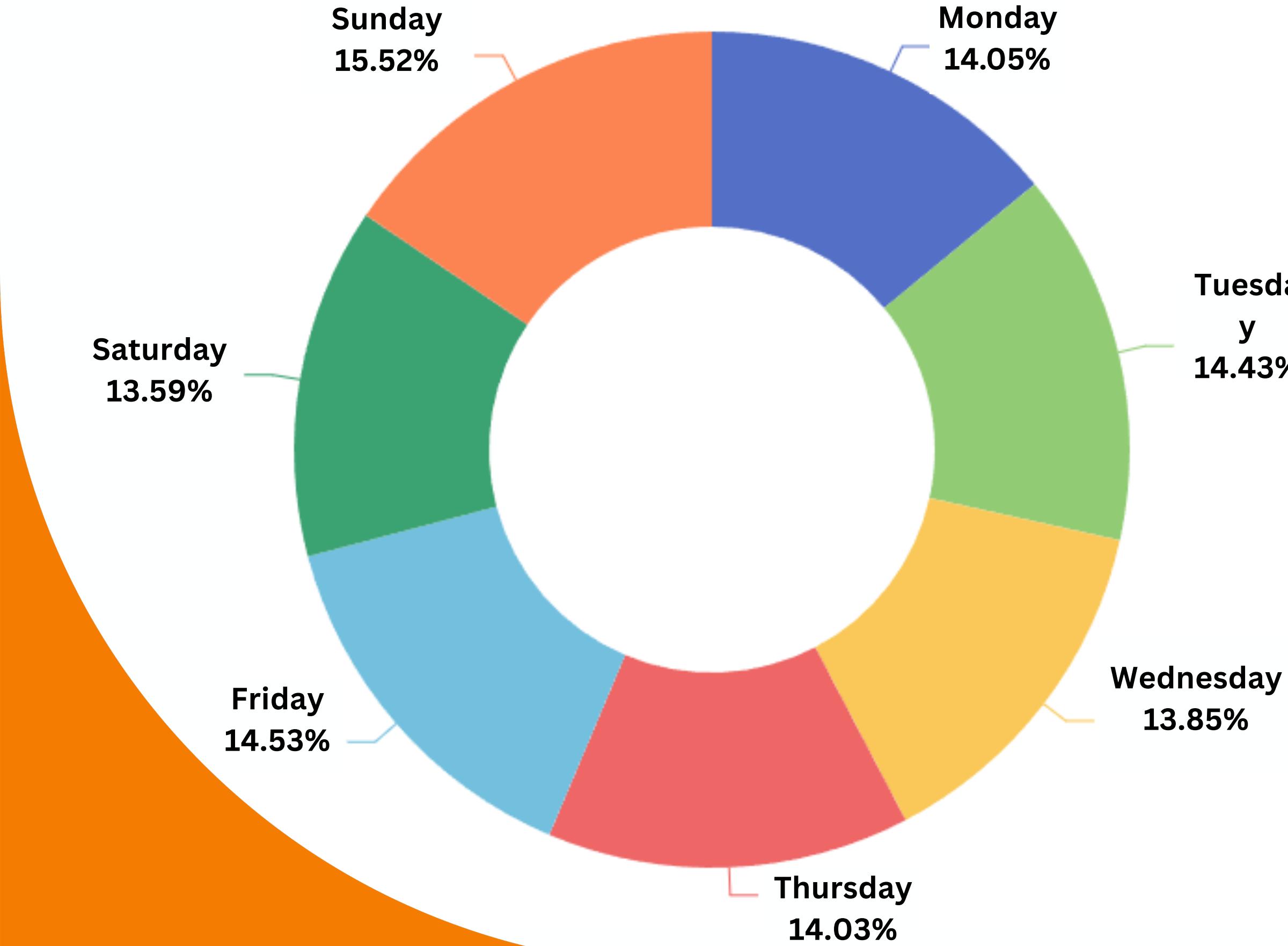
Observations:

- The number of messages decreased steadily by 90% in the past 3 years
- less chat at the beginning of a semester instead of “blood final”
- more chat during quarantine

### 3. Pattern Analysis - Week



Distribution of Messages in a Week



- My expectation:  
Peak during weekends

**V.S.**

- Reality  
Uncorrelated  
with weekends

# 3. Pattern Analysis - Day



## Peak Hour in a Day

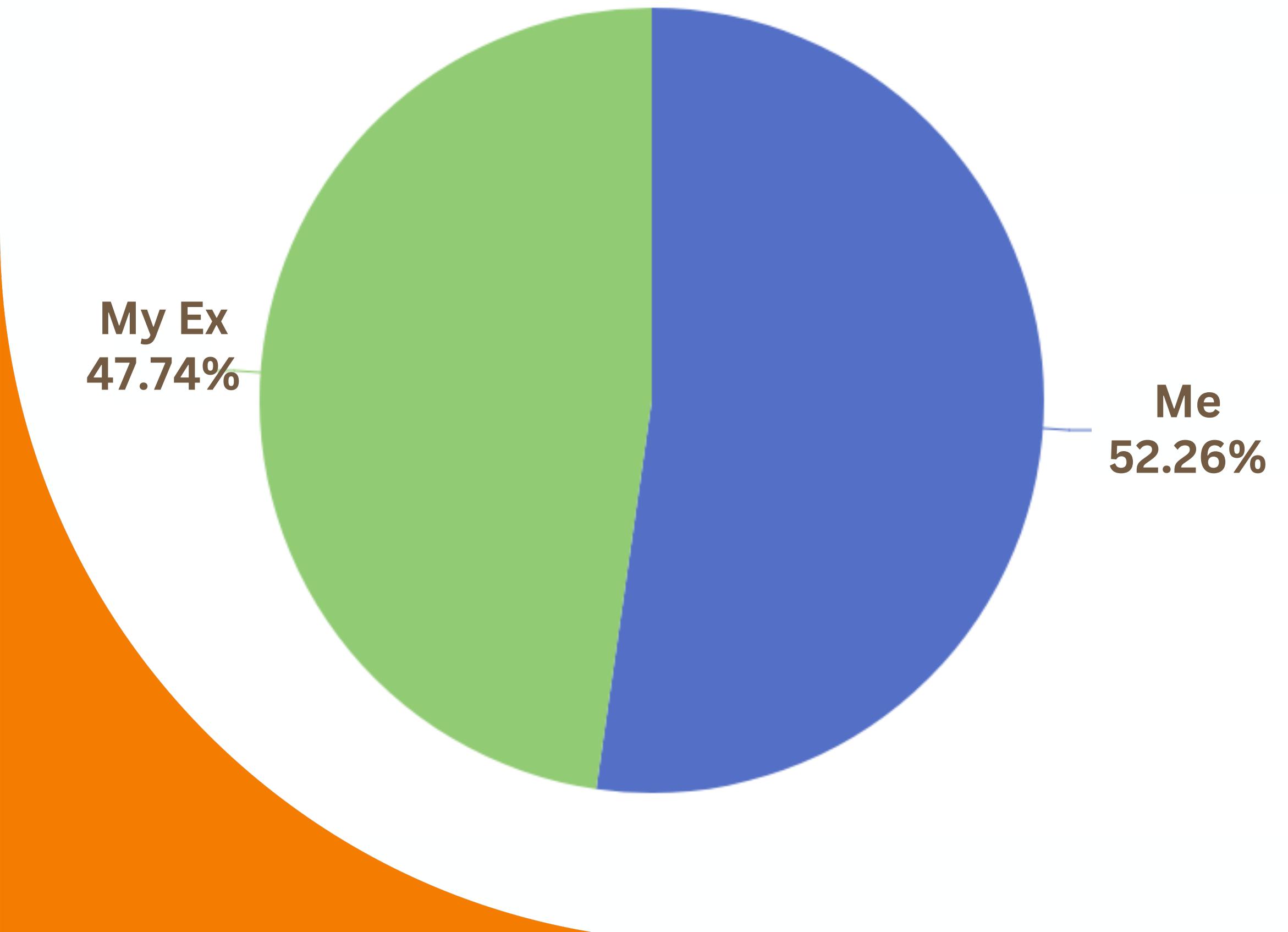


Chat happens most  
during lunch time  
and before sleeping

### 3. Pattern Analysis - Sender



Percentage of Messages from Each Sender



- My expectation:  
 $Me = 3 * \text{Him}$

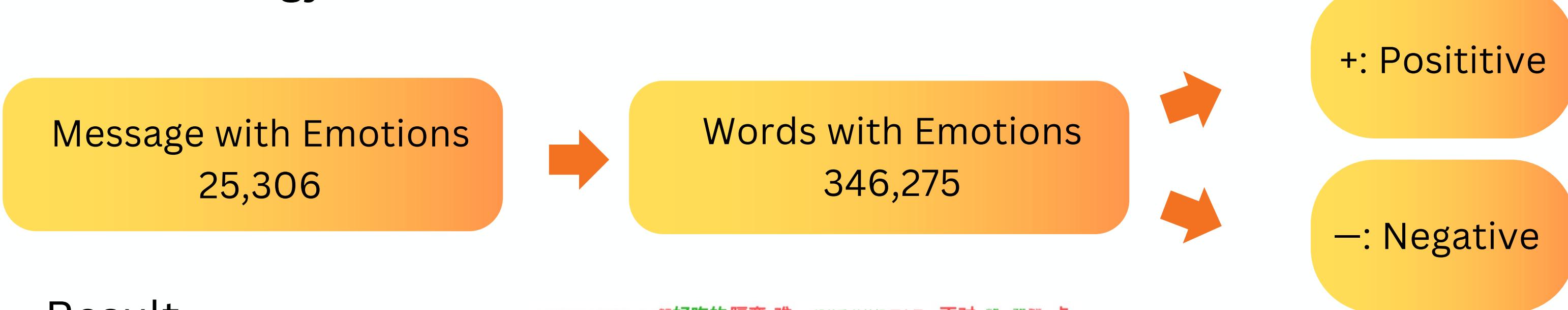
**v.s.**

- Reality  
 $Me \approx \text{Him}$

# 4. Sentiment Analysis



- Customize the sentiment value of frequently used words
  - From negative to positive:
    - 好傻/傻傻 (silly)
    - 肥肥鲨 (fat shark)
- Methodology

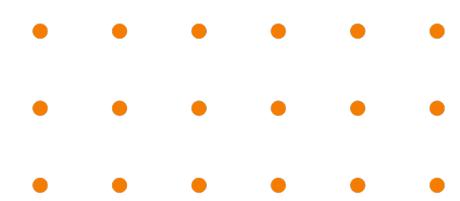


- Result

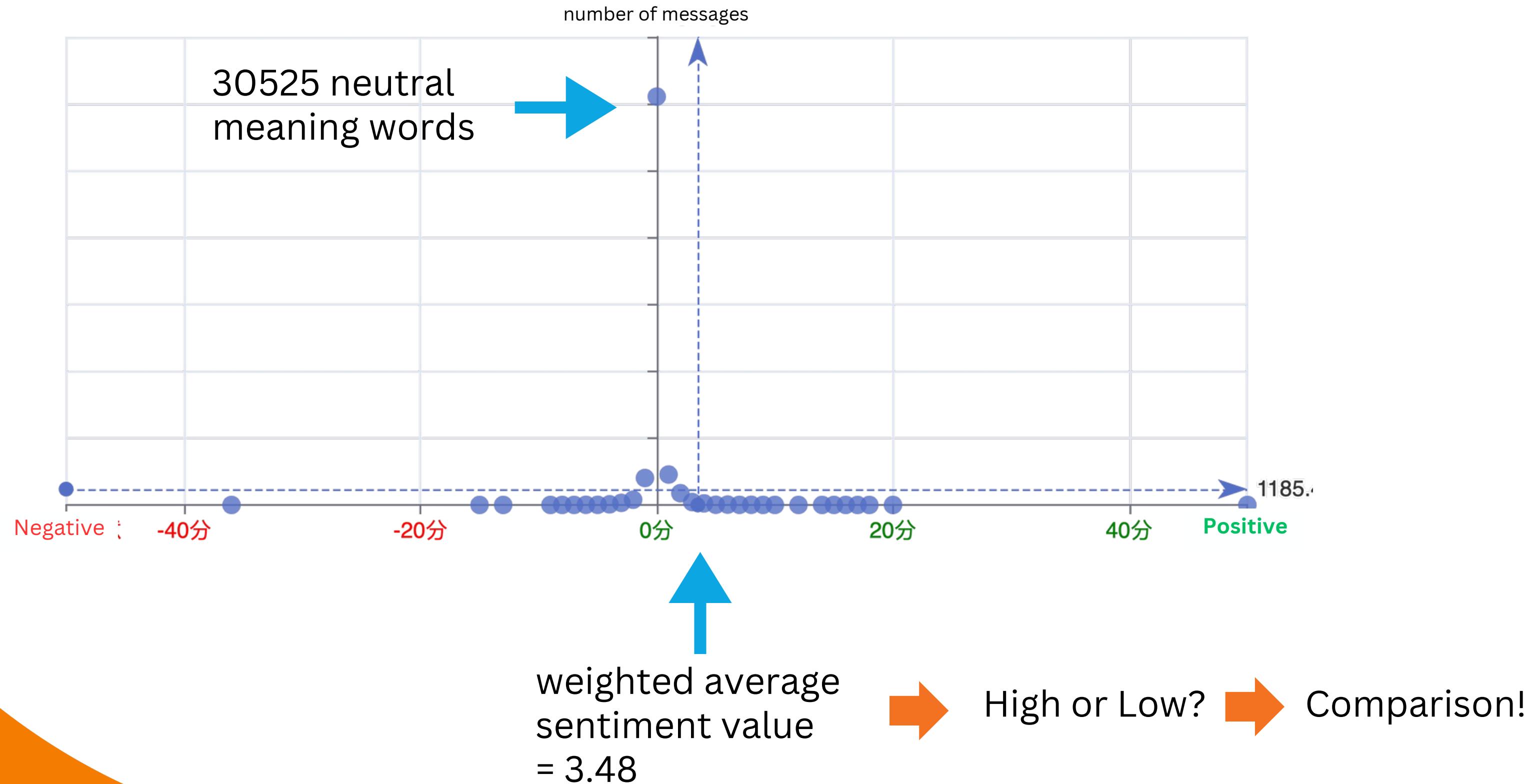
感染 喜欢的 涩<sup>辛微</sup> 好吃的 隔离 难<sup>明白</sup> 没关系 比较好 不方便 不对<sup>当然</sup> 了解 取消 卡<sup>标准</sup>  
正常 大哭 生气 具体 疫情 哎 好累 哈 不舒服 辛苦 感冒 奇怪 不懂  
哼! 鲨鱼<sup>麻烦</sup> 厉害 很多 爱 害羞 怪 微笑 学习 很快 难受 烟花 嘲讽 打开  
注意<sup>欢迎</sup> 好多 没到<sup>清楚</sup> 不行 完了 裂开 主要 最好<sup>严重</sup> 发现 确实 完全 活动<sup>自然</sup>  
机会<sup>重惊</sup> 你妈<sup>朋友</sup> 不太<sup>为什么</sup> 不知道<sup>不知道</sup> 休息 好吃 好玩 痛<sup>痛</sup>  
少<sup>舒服</sup> 适合<sup>肯定</sup> 不太好<sup>不太好</sup> 宝贝 加油 哈哈 忘记<sup>解决</sup> 乐<sup>错过</sup> 可爱的<sup>满足</sup> 甜  
暂时<sup>敲打</sup> 快乐 相信 刚刚 可爱<sup>怕</sup> 接受 爱你 喜欢 需要 随便<sup>顺利</sup>  
便宜<sup>真棒</sup> 去你<sup>捂脸</sup> 其实 没有<sup>不是</sup> 不可能 一定 坏 不给 可怜<sup>老是</sup>  
参加<sup>不同</sup> 不好<sup>诱惑</sup> 不是<sup>可能</sup> 一定 坏 不给 可怜<sup>老是</sup>  
不许<sup>不帮</sup> 不错<sup>不理</sup> 没事<sup>不能</sup> 未知<sup>准备</sup> 知道<sup>事情</sup> 好奇<sup>控制</sup> 还好<sup>确定</sup> 刚好<sup>刚好</sup>  
辛苦了<sup>流泪</sup> 好好<sup>苦涩</sup> 傻 诶<sup>问题</sup> 不会<sup>不会</sup> 傻瓜 算了<sup>一口</sup> 差<sup>容易</sup> 不好意思<sup>更好</sup>  
寂寞<sup>好不容易</sup> 咖啡<sup>嘿嘿</sup> 棒<sup>打算</sup> 一般<sup>明郎</sup> 不要<sup>拒绝</sup> 不了<sup>方便</sup> 好看 开心 担心 才能<sup>玫瑰</sup>  
拜拜<sup>绝对</sup> 特别<sup>理解</sup> 疯狂<sup>突然</sup> 希望<sup>没想到</sup> 超级<sup>超级</sup> 还行<sup>发烧</sup> 老师<sup>贵</sup> 谢谢<sup>绝了双气</sup>  
帮忙<sup>专门</sup> 考虑<sup>不喜欢</sup> 笑死<sup>挺好的</sup> 推<sup>说明</sup> 期待<sup>讨厌</sup> 暗示<sup>清瘦</sup> 睡不着<sup>发烧</sup> 好好看<sup>通知</sup> 看到家<sup>有啥</sup> 庆祝<sup>快速</sup>



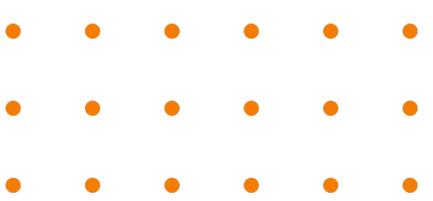
# 4. Sentiment Analysis



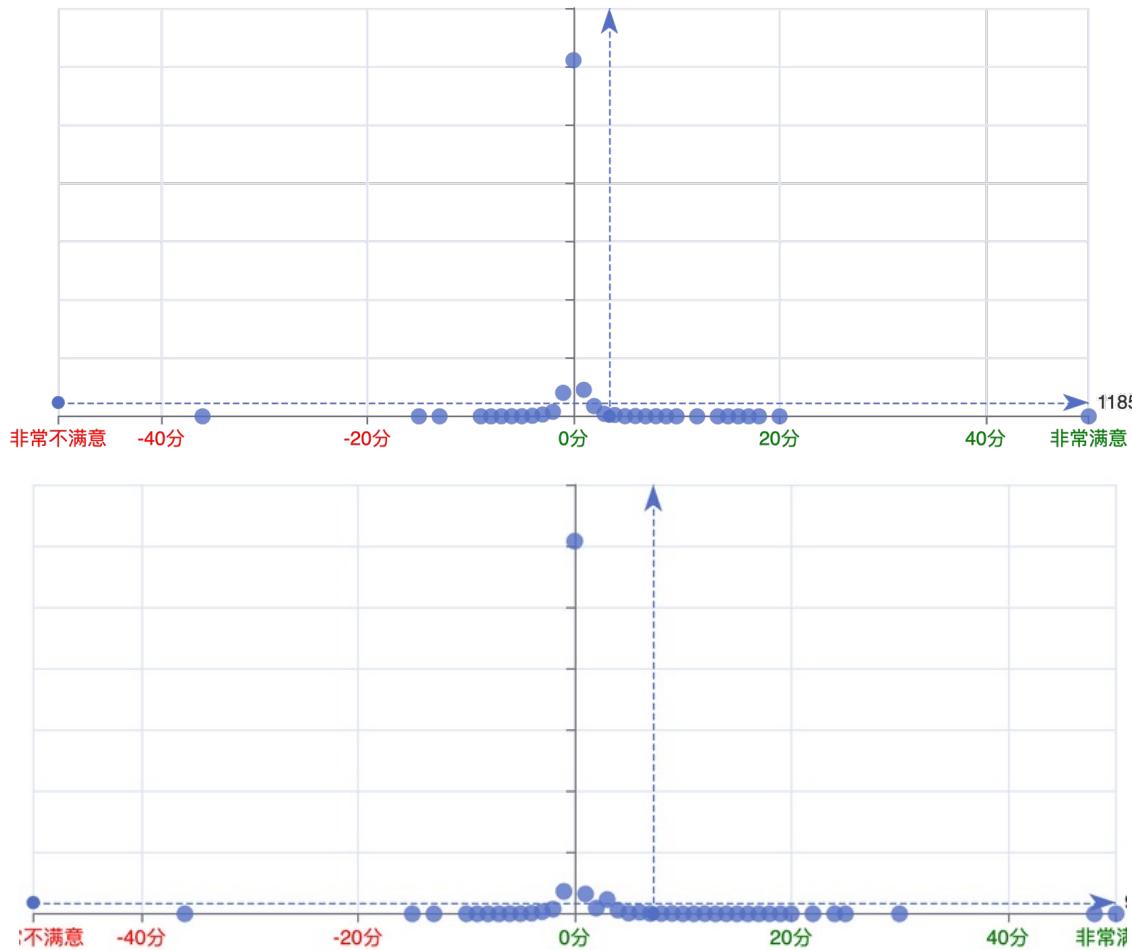
## Distribution of Words with Different Sentiment Values



# 4. Sentiment Analysis



Comparison of Weighted Average Sentiment Value



Data	Time	Value
Me & my Ex	First 8 years	3.48
My friend & her Boyfriend	First 3 years	7.2
Me & my Ex	First 3 years	?

Comparison Group

- Long-Distance Relationship too
- Business & Medical Students too
- Similar personality

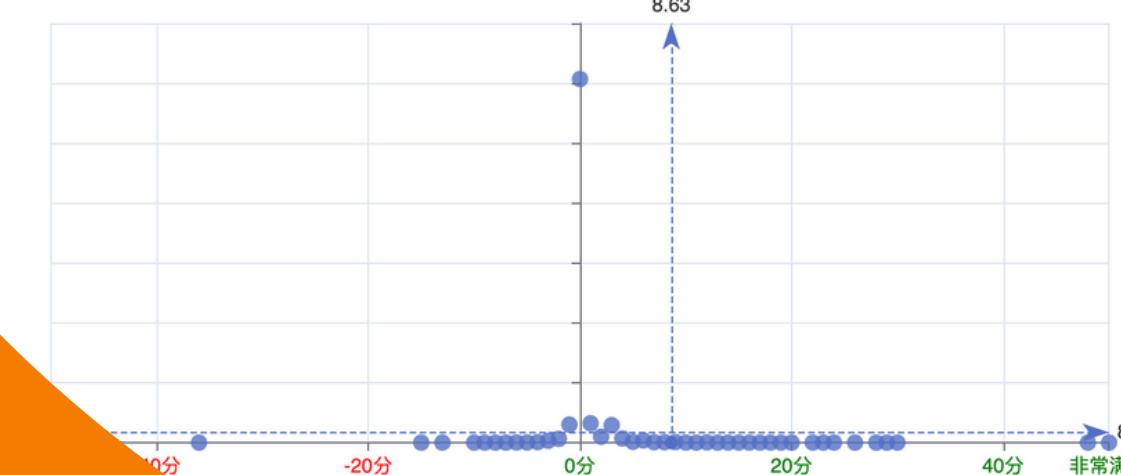
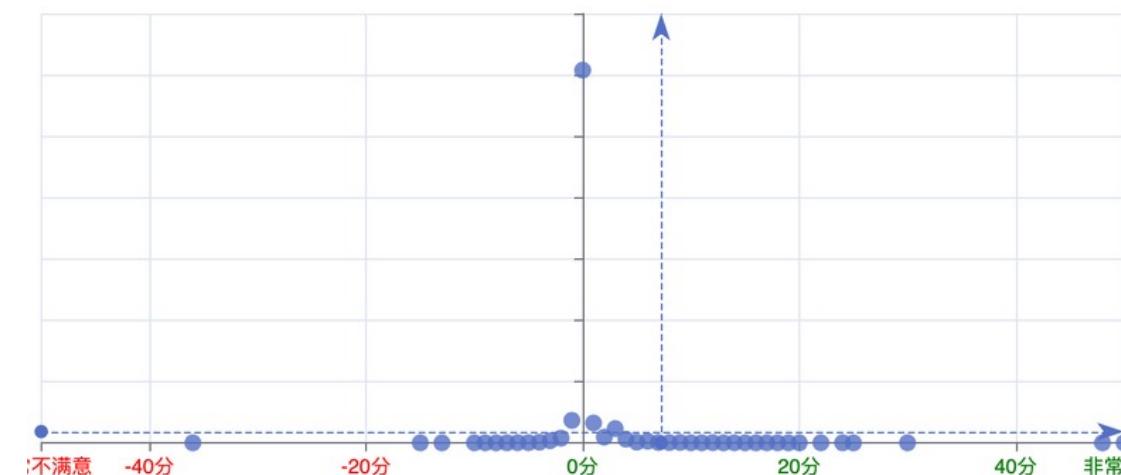
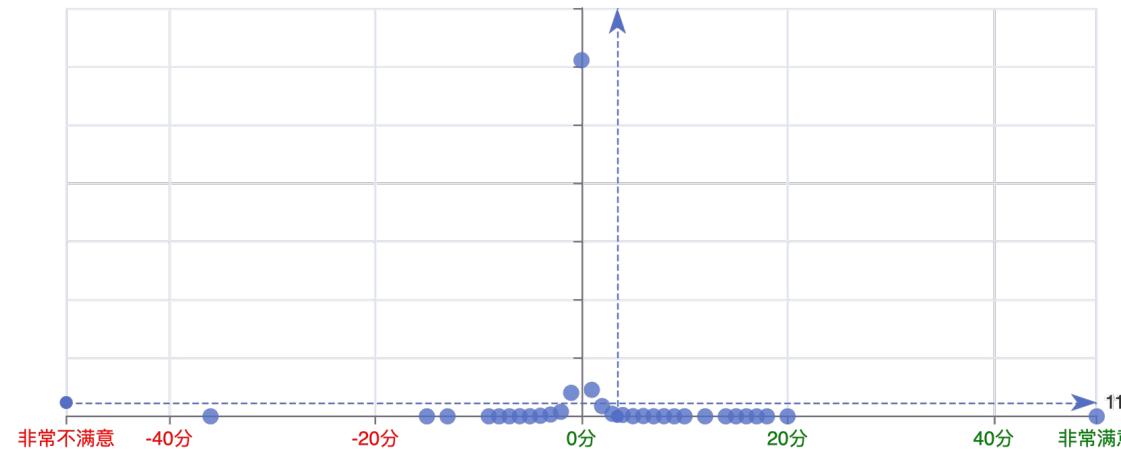
MBTI:

ENFJ ♀ & INTP ♂

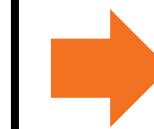
# 4. Sentiment Analysis



Comparison of Weighted Average Sentiment Value

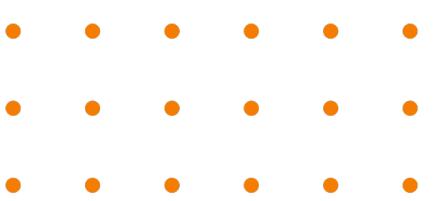


Data	Time	Value
Me & my Ex	First 8 years	3.48
My friend & her Boyfriend	First 3 years	7.2
Me & my Ex	First 3 years	8.63

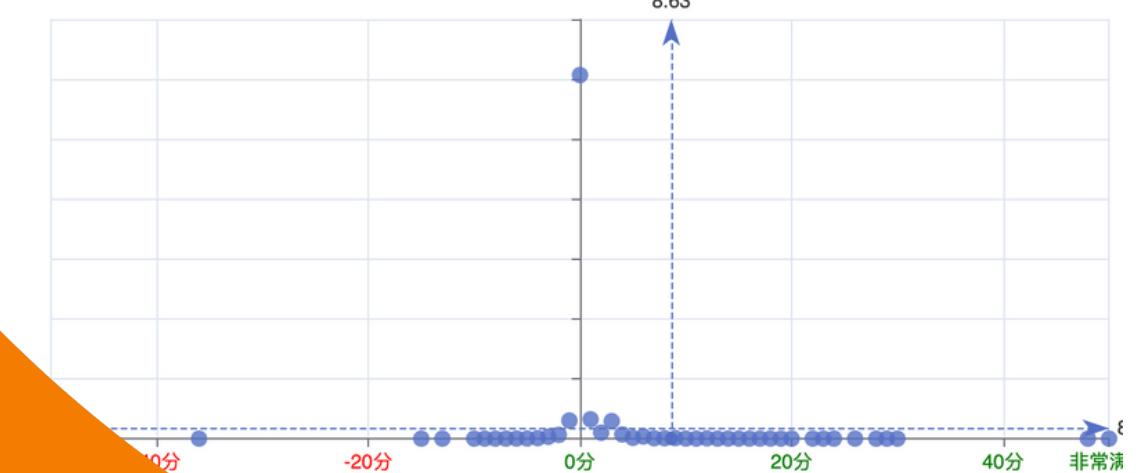
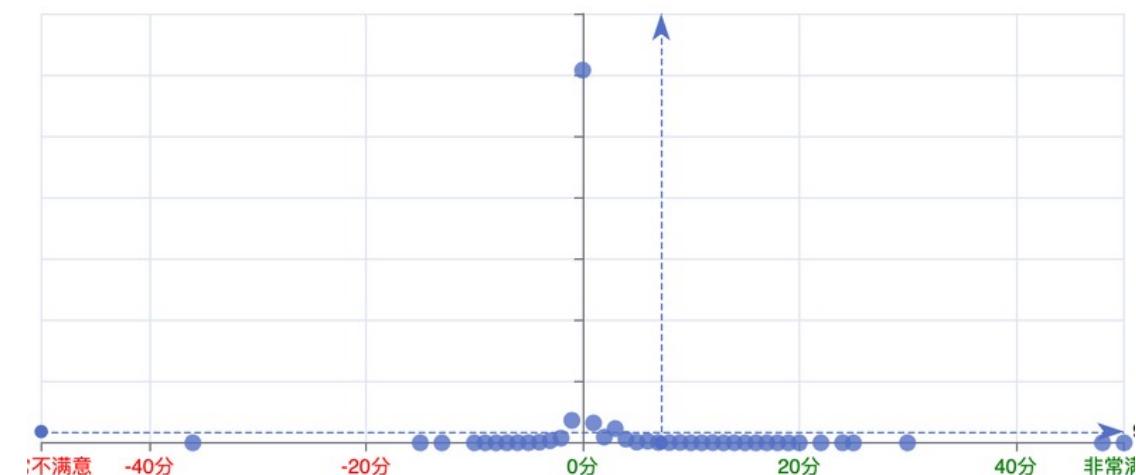
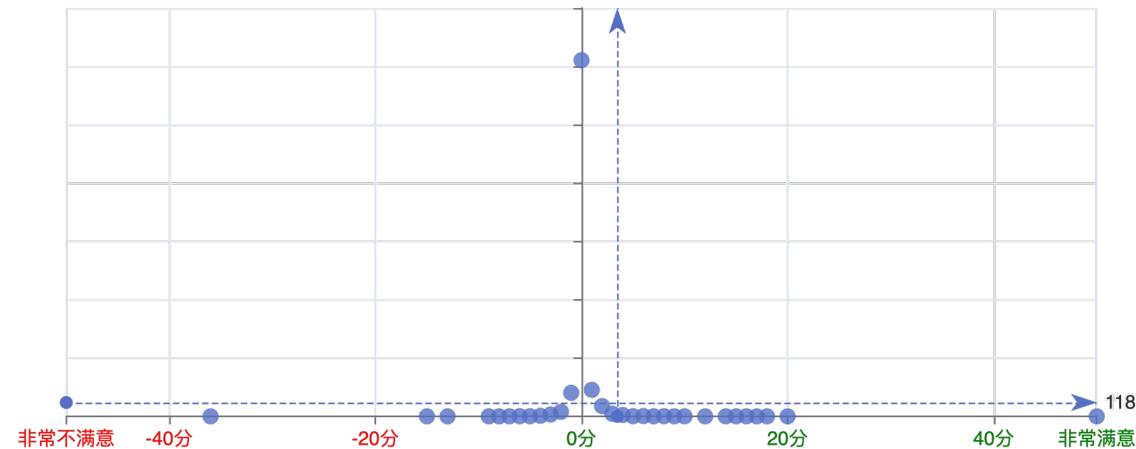


After controlling the *Time*, our data outperforms

# 4. Sentiment Analysis



Comparison of Weighted Average Sentiment Value



Data	Time	Value
Me & my Ex	First 8 years	3.48
My friend & her Boyfriend	First 3 years	7.2
Me & my Ex	First 3 years	8.63



After controlling the *Time*, our data outperforms

*Time kills Everything* 😂

# 5. Conclusion – What Kills Our Relationship

- **Busyness**

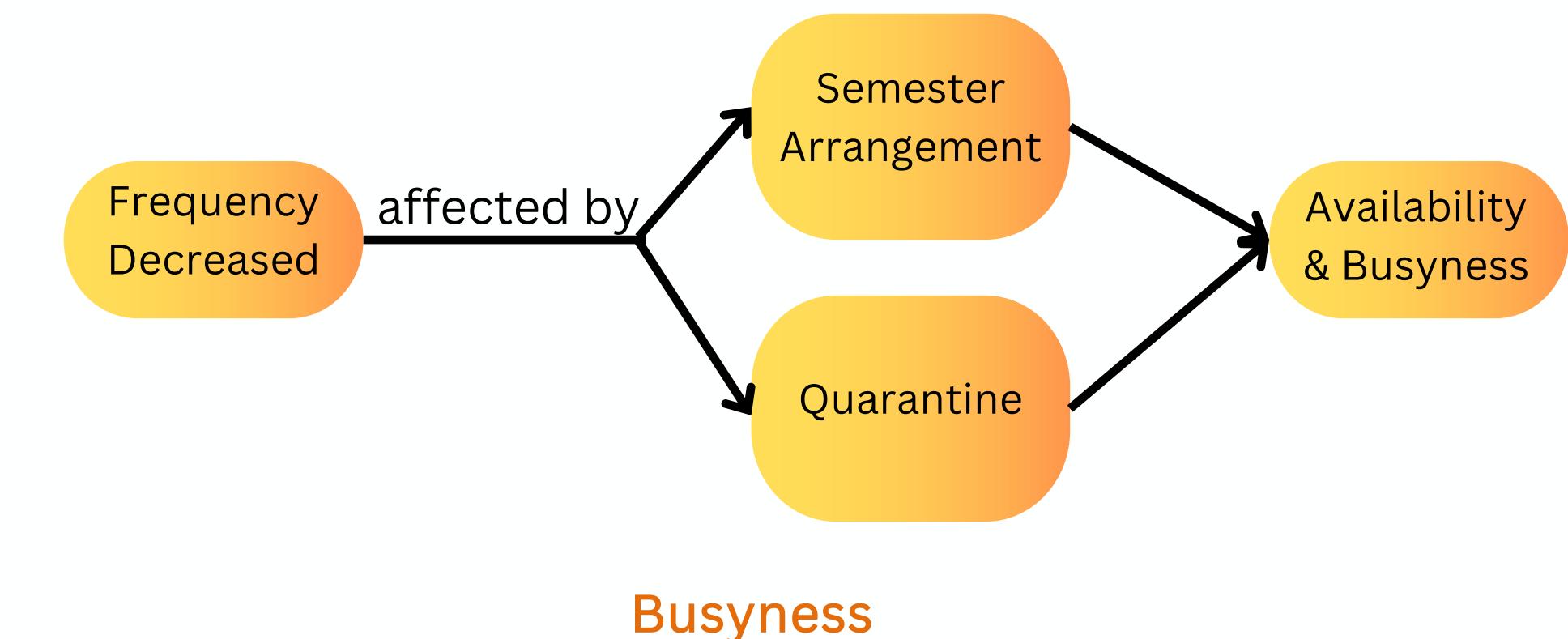
Time availability is the top factor affecting our chat frequency

- **Misperception**

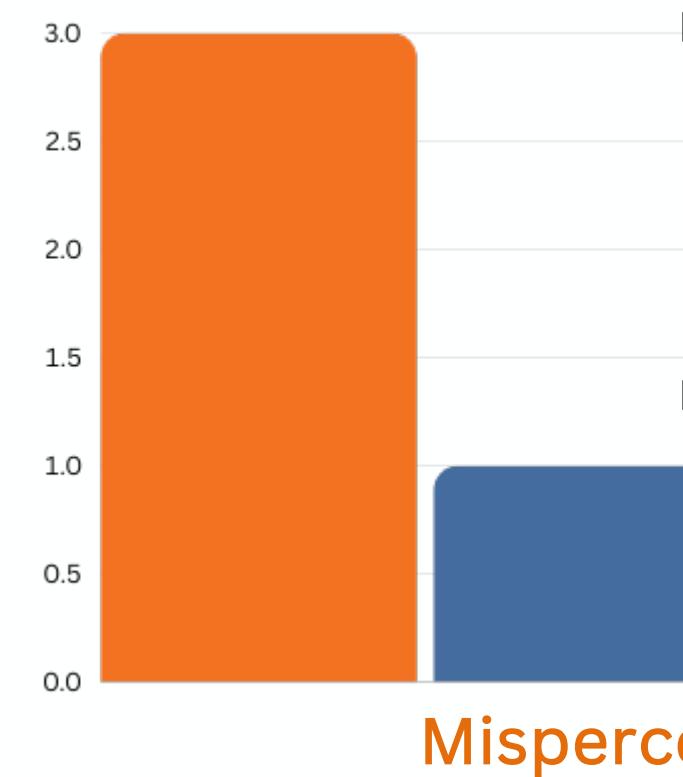
The emotional delusion leads to deviation of perception from reality

- **Power of Time**

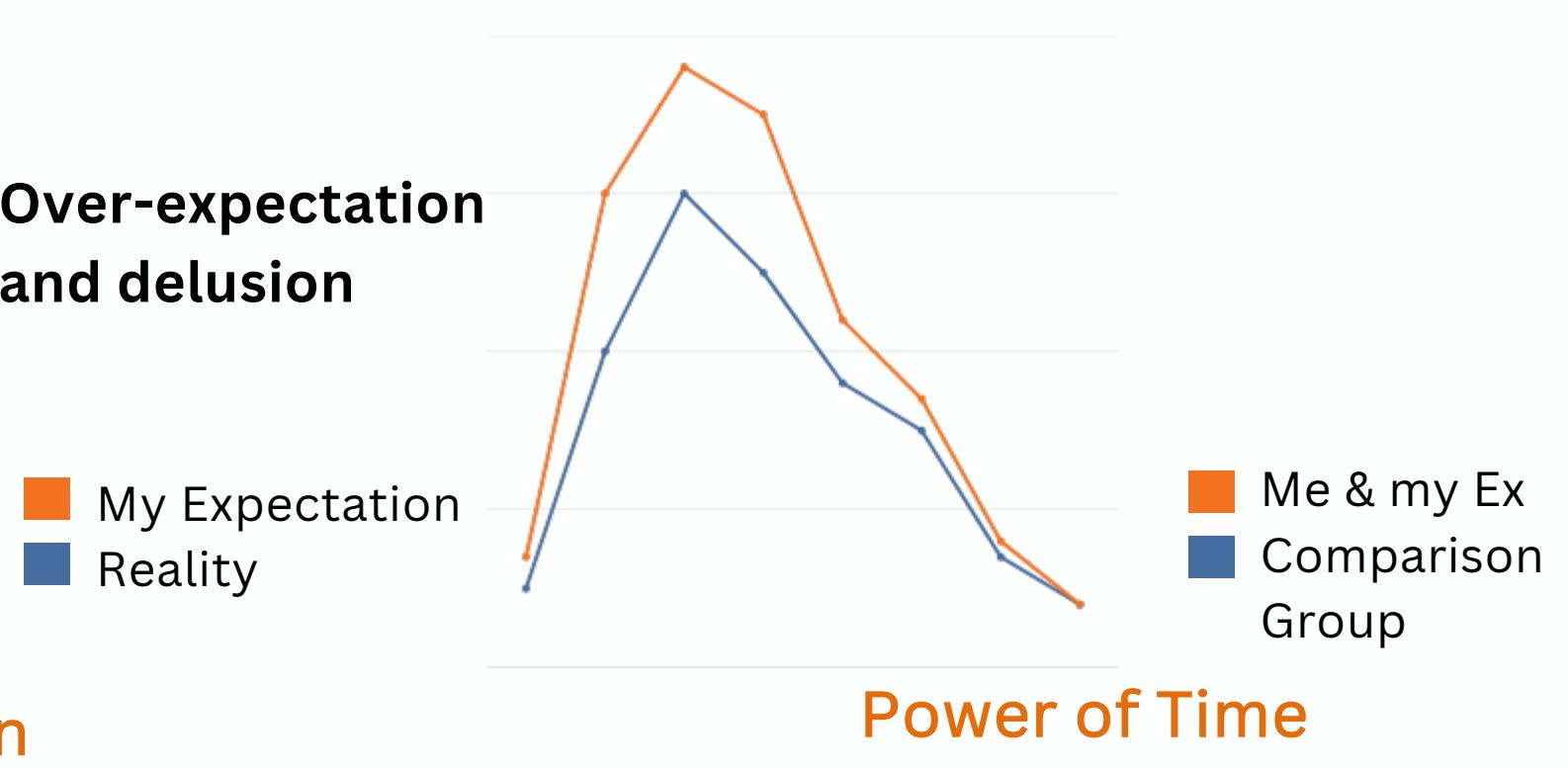
Unavoidable fading of passion



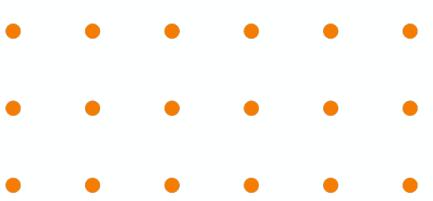
Am I More Passionate than Him?



Evolution of Passion



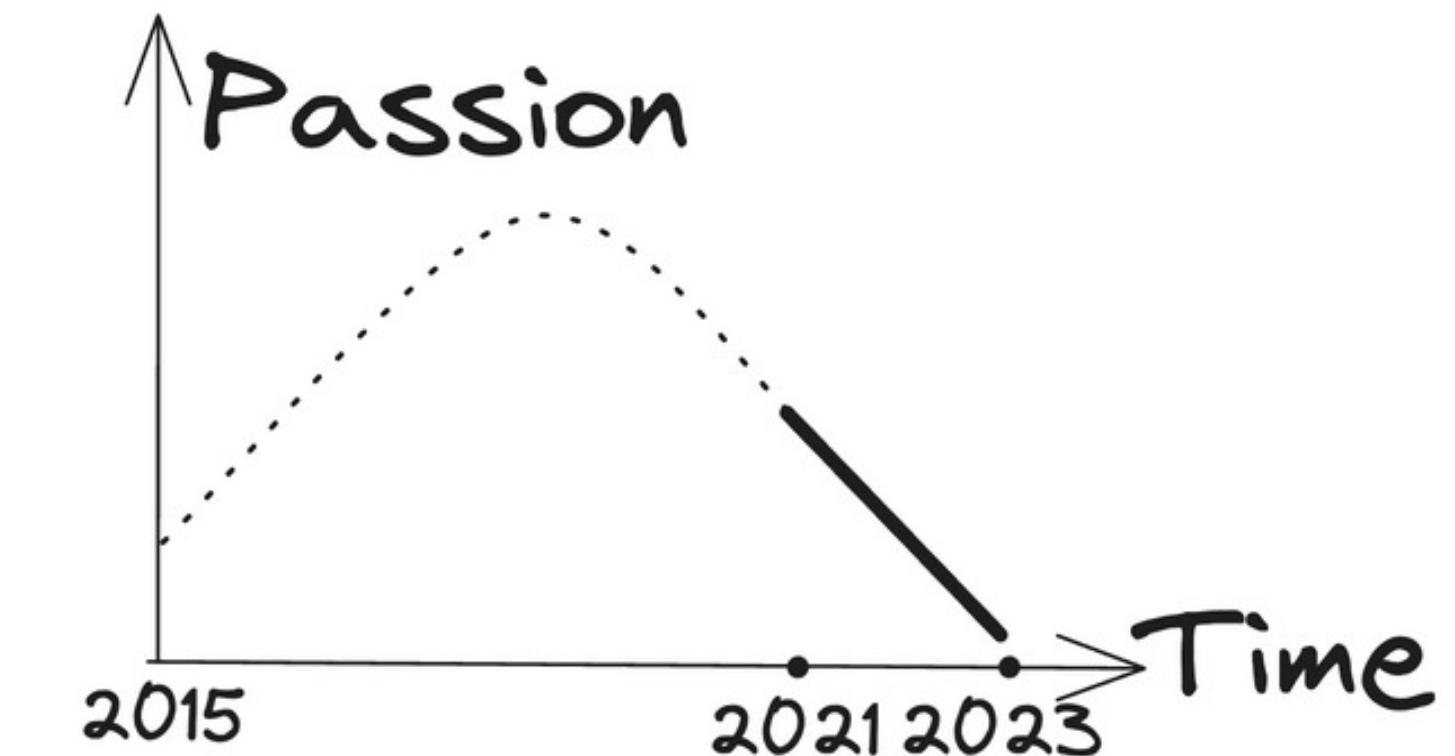
# 6. Limitation



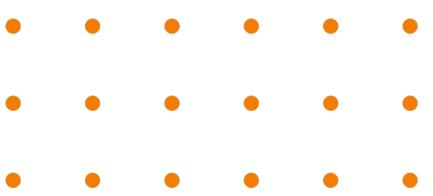
- Misunderstanding of Language Reduces the Accuracy in understanding the meaning

Type	Example	Wrongly Classified as
Sarcasm	It's <b>OK</b> (after getting stood up)	Positive

- Lack of the whole-life data
  - The information before 2021 May is absent, which disables the whole-life analysis on this 8-year-old relationship, as well as the comparison before and after the turning point

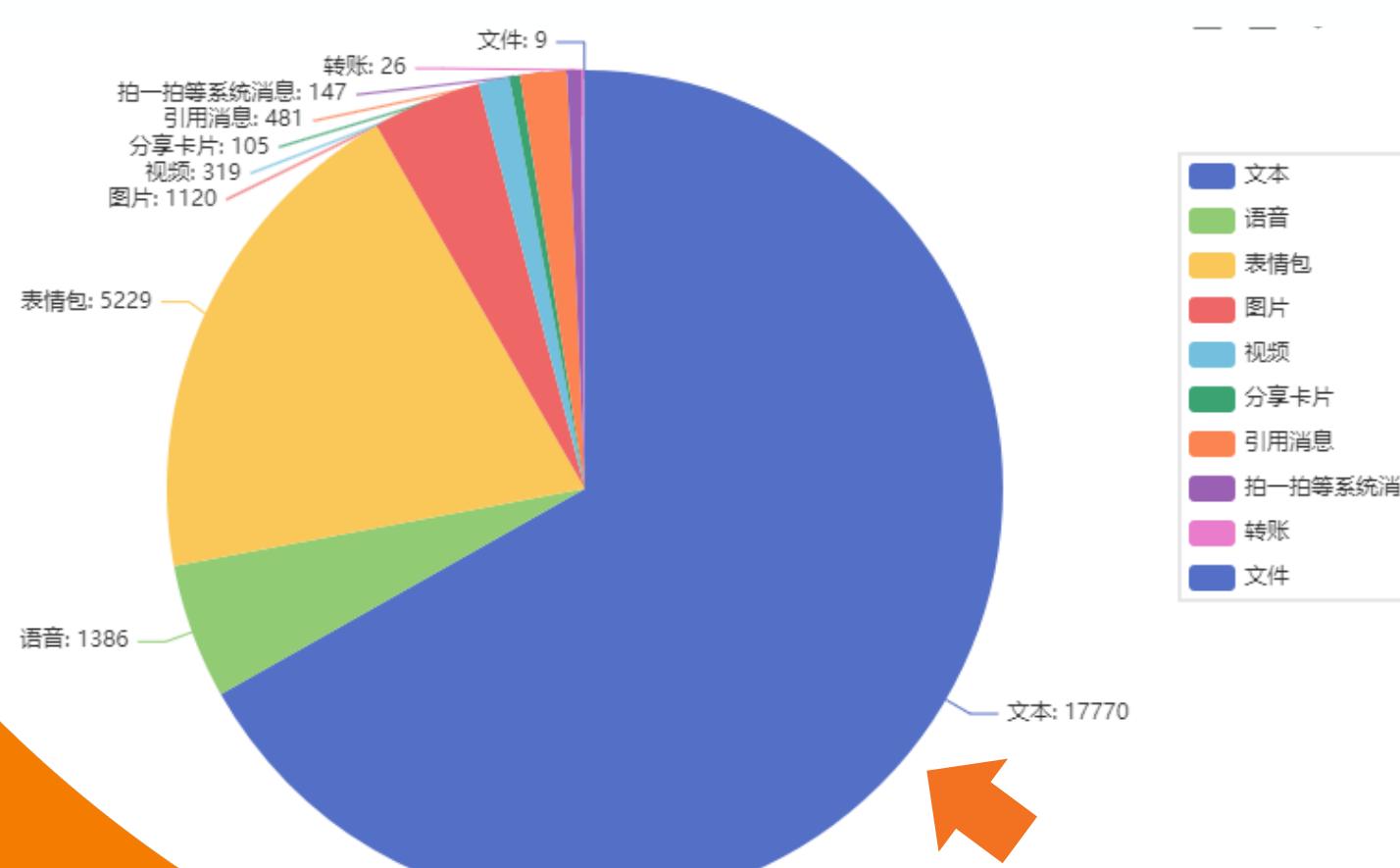


# 6. Limitation



- Bias of Partial Selection
  - Only string data is included (67.69%). Many emotion-bearing data types are ignored
  - Example:
    - Frequent “container” of emotion: emoji
    - On significant scenarios, voice messages or phone calls are preferred (e.g. on birthday, after receiving HKU offer)

**Text only accounts for 67.79%**



**Most Used Emoji: Cuddle**



Cuddle

# Thank you!



Presented By:  
Cheng Ling Jun