

# Dictionary в Python

По мотивам Objects/dictnotes.txt

Cyril @notorca Lashkevich

# Как создать словарь

- {}
- dict()
- PyObject\* PyDict\_New()

# Сколько словарей в Hello World?

```
$ python -c "print('Hello world')"  
| wc -l
```

# Сколько словарей в Hello World?

- \$ python -c "print('Hello world')"  
| wc -l
- 1642

# Именованные параметры функций

- 1 запись, 1 чтение
- 1-3 элемента
- Часто встречается в обычных программах на Python

# Поиск метода в классе

- 1 запись, много чтений
- 8-16 элементов
- При наследовании много неудачных чтений с последующим поиском в базовом классе

# Атрибуты и глобальные пременные

- Много записей и чтений
- 4-10 элементов

# Builtins

- Частые чтение, почти не бывает записи
- ~150 строковых ключей (3.3)
- По некоторым ключам чтения гораздо чаще чем по другим

# Удаление повторов, подсчет элементов

- Однократное чтение по каждому из ключей
- Произвольное количество элементов
- Многократный доступ по одному ключу подряд

# Удаление дубликатов

`dict.fromkeys(seqn).keys()`

- Все операции записи при конструировании

# Подсчет элементов в последовательности

```
for e in seqn:  
    d[e] = d.get(e,0) + 1
```

- 2 последовательных доступа по одинаковому ключу

# Создание индекса из словаря списков

```
for pnum, page in enumerate(pages):  
    for w in page:  
        d.setdefault(w, []).append(pnum)
```

- `setdefault` совмещает 2 поиска в 1м

# Проверка принадлежности

- Словари произвольных размеров
- Создаются 1 раз и затем мало изменяются
- Много вызовов `has_key()` и  
`__contains__()`

# Динамические отображения

- Чередующиеся добавления, удаления, чтение и перезапись элементов

# Реализация (2.7)

- Последовательная область памяти с доступом по индексу

```
typedef struct {  
    Py_hash_t me_hash;  
    PyObject *me_key;  
    PyObject *me_value;  
} PyDictKeyEntry;
```

# Пустой dict с размером по умолчанию (8 элементов)

```
>>> d = {}
```

Idx	Hash	Key	Value
000			
001			
010			
011			
100			
101			
110			
111			

# Хеширование ключа

- Ключ преобразуется в индекс с помощью функции `hash()`
- `hash()` возвращает 32/64bit значение
- Для индекса берется n младших бит

# Свойства хеша

- Для равных значений хеши всегда равны
- Даже если представление значений разное: 9, 9.0, complex(9,0)
- Похожие значения дают сильно отличающиеся хеши

```
>>> d['ftp'] = 21  
>>> bits(hash('ftp'))[-8:]  
10100001
```

Idx	Hash	Key	Value
000			
001			
010			
011			
100			
101			
110			
111			

```
>>> d['ftp'] = 21
```

```
>>> bits(hash('ftp'))[-8:]
```

10100001

Idx	Hash	Key	Value
000			
001	= ...10100001	'ftp'	21
010			
011			
100			
101			
110			
111			

```
>>> d['ssh'] = 22
```

```
>>> bits(hash('ssh'))[-3:]
```

```
101
```

Idx	Hash	Key	Value
000			
001	= ...10100001	'ftp'	21
010			
011			
100			
101			
110			
111			



```
>>> d['ssh'] = 22
>>> bits(hash('ssh'))[-3:]
101
```

Idx	Hash	Key	Value
000			
001	= ...10100001	'ftp'	21
010			
011			
100			
→ ● 101	= ...01010101	'ssh'	22
110			
111			

```
>>> d['smtp'] = 25  
>>> bits(hash('smtp'))[-3:]  
100
```

Idx	Hash	Key	Value
000			
001	= ...10100001	'ftp'	21
010			
011			
100			
101	= ...01010101	'ssh'	22
110			
111			



```
>>> d['smtp'] = 25  
>>> bits(hash('smtp'))[-3:]  
100
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010				
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111				

```
>>> d['time'] = 37  
>>> bits(hash('time'))[-3:]  
111
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010				
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111				



```
>>> d['time'] = 37  
>>> bits(hash('time'))[-3:]  
111
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010				
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111	=	...11110111	'time'	37

```
>>> d['www'] = 80  
>>> bits(hash('www'))[-3:]  
010
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010				
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111	=	...11110111	'time'	37



```
>>> d['www'] = 80  
>>> bits(hash('www'))[-3:]  
010
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010	=	...10000010	'www'	80
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111	=	...11110111	'time'	37



```
d = {'ftp': 21, 'ssh': 22,  
      'smtp': 25, 'time': 37,  
      'www': 80}
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010	=	...10000010	'www'	80
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111	=	...11110111	'time'	37

# Поиск в словаре

- Вычислить хеш от ключа
- Обрезать старшие биты
- Взять значение из слота по индексу

```
>>> d['smtp']
25
>>> bits(hash('smtp'))[-3:]
100
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010	=	...10000010	'www'	80
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111	=	...11110111	'time'	37

# Перебор всех элементов

- Словари возвращают свои ключи или значения в порядке отличном от порядка добавления

```
>>> print d  
{'ftp': 21, 'www': 80, 'smtp':  
25, 'ssh': 22, 'time': 37}
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010	=	...10000010	'www'	80
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111	=	...11110111	'time'	37

```
>>> d.keys()  
['ftp', 'www', 'smtp', 'ssh',  
'time']
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010	=	...10000010	'www'	80
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111	=	...11110111	'time'	37

```
>>> d.values()  
[21, 80, 25, 22, 37]
```

Idx		Hash	Key	Value
000				
001	=	...10100001	'ftp'	21
010	=	...10000010	'www'	80
011				
100	=	...10001100	'smtp'	25
101	=	...01010101	'ssh'	22
110				
111	=	...11110111	'time'	37

# Коллизии

- Разные ключи пытаются доступиться по одинаковому индексу
- Находим первое свободное место

```
>>> d = {}
```

Idx	Hash	Key	Value
000			
001			
010			
011			
100			
101			
110			
111			

```
>>> d['smtp'] = 21
```

Idx	Hash	Key	Value
000			
001			
010			
011			
100			
101			
110			
111			

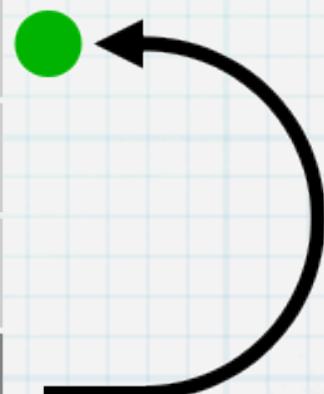


```
>>> d['smtp'] = 21
```

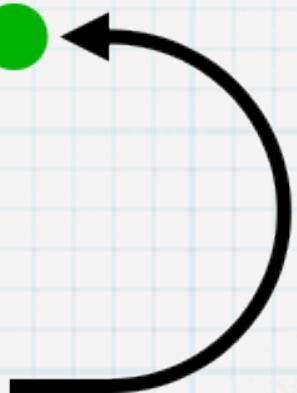
Idx	Hash	Key	Value
000			
001			
010			
011			
100	= ...10001100	'smtp'	21
101			
110			
111			

```
>>> d['dict'] = 2628
```

Idx	Hash	Key	Value
000			
001			
010			
011			
100	= ...10001100	'smtp'	21
101			
110			
111			



```
>>> d['dict'] = 2628
```

Idx		Hash	Key	Value	
000					
001	≠	...00110100	'dict'	2628	
010					
011					
100	→ x	= ...10001100	'smtp'	21	
101					
110					
111					

```
>>> d['svn'] = 3690
```

Idx	Hash	Key	Value
000			
001	≠ ...00110100	'dict'	2628
010			
011			
100	= ...10001100	'smtp'	21
101			
110			
111			

```
>>> d['svn'] = 3690
```

	Idx	Hash	Key	Value
→ ●	000	= ...11100000	'svn'	3690
	001	≠ ...00110100	'dict'	2628
	010			
	011			
	100	= ...10001100	'smtp'	21
	101			
	110			
	111			

```
>>> d['ircd'] = 6667
```

Idx	Hash	Key	Value
000	= ...11100000	'svn'	3690
001	≠ ...00110100	'dict'	2628
010			
011			
100	= ...10001100	'smtp'	21
101			
110			
111			

Diagram illustrating the search and insertion process:

- An arrow points to index 000, where the value 3690 is highlighted.
- A red circle with an 'x' is placed over the '000' index.
- The '000' index cell contains a green '=' symbol.
- The '001' index cell contains a red '≠' symbol.
- Three black arrows point from the 'x' at index 000 to the 'x' at index 100, indicating a collision.
- A green circle is placed at index 111, indicating it is available for insertion.

```
>>> d['ircd'] = 6667
```

Idx	Hash	Key	Value
000	= ...11100000	'svn'	3690
001	≠ ...00110100	'dict'	2628
010			
011			
100	= ...10001100	'smtp'	21
101			
110	≠ ...01001000	'ircd'	6667
111			

Diagram illustrating the search and insertion process:

- An arrow points to index 000, where the value 3690 is crossed out (red 'x') and replaced by 6667 (green circle).
- Three red 'x' marks are placed on the row for index 001.
- Three black curved arrows point from the bottom right towards the row for index 110, indicating the search path for the new key 'ircd'.

```
>>> d['zope'] = 9673
```

Idx	Hash	Key	Value
000	= ...11100000	'svn'	3690
001	≠ ...00110100	'dict'	2628
010			
011			
100	= ...10001100	'smtp'	21
101			
110	≠ ...01001000	'ircd'	6667
111			

Diagram illustrating the search and insertion process:

- An arrow points to index 000, where the value 3690 is highlighted.
- A red circle with an 'x' is placed over the '000' index.
- A green circle highlights the '=' sign in the Hash column for index 000.
- A red circle with an 'x' is placed over the '001' index.
- A green circle highlights the '2628' value in the Value column for index 001.
- A black curved arrow starts from the '2628' value and loops back to the '001' index.

```
>>> d['zope'] = 9673
```

# 2 из 5ти элементов на своих  
ожидаемых местах

Idx	Hash	Key	Value
000	...11100000	'svn'	3690
001	...00110100	'dict'	2628
010			
011	...10111000	'zope'	9673
100	...10001100	'smtp'	21
101			
110	...01001000	'ircd'	6667
111			

Diagram annotations:

- A red arrow points to the first row (Idx 000). A red circle with an 'x' is placed over the 'Idx 000' cell.
- A green square with an '=' sign is placed over the 'Idx 000' cell.
- A red circle with an 'x' is placed over the 'Idx 011' cell.
- A green circle is placed over the 'Idx 011' cell.
- A black curved arrow starts from the 'Value' column of the 'zope' row and loops back to the 'Idx' column of the same row.

# Коллизии и очередьность

- Поскольку из-за коллизий элементы могут находиться не по своим "естественным" индексам порядок элементов зависит от порядка добавления

# Поиск первой свободной ячейки

- Последовательный поиск плох для int ключей

```
pertrurb = hash  
while (<слот занят>) {  
    slot = (5*slot) + 1 + perturb;  
    perturb >>= 5;  
}
```

```
>>> d['svn']
```

3690

	Idx	Hash	Key	Value
→ ●	000	= ...11100000	'svn'	3690
	001	≠ ...00110100	'dict'	2628
	010			
	011	≠ ...10111000	'zope'	9673
	100	= ...10001100	'smtp'	21
	101			
	110	≠ ...01001000	'ircd'	6667
	111			

```
>>> d['ircd']
```

```
6667
```

Idx	Hash	Key	Value
000	= ...11100000	'svn'	3690
001	≠ ...00110100	'dict'	2628
010			
011	≠ ...10111000	'zope'	9673
100	= ...10001100	'smtp'	21
101			
110	≠ ...01001000	'ircd'	6667
111			

Diagram illustrating the search process:

- An arrow points to index 000, which has a red circle with an 'x' over it.
- The cell at index 001 has a red circle with an 'x' over it.
- The cell at index 100 has a red circle with an 'x' over it.
- The cell at index 110 has a green circle with an 'x' over it.
- A green circle is placed next to the value 6667 at index 110.
- Curved arrows point from the 'x' marks to the right edge of the table, indicating they are being discarded.

```
>>> d['nsca']  
KeyError: 'nsca'
```

Idx		Hash	Key	Value
000	=	...11100000	'svn'	3690
001	≠	...00110100	'dict'	2628
010				
011	≠	...10111000	'zope'	9673
100	=	...10001100	'smtp'	21
101				
110	≠	...01001000	'ircd'	6667
111				



```
>>> d['netstat']
```

```
KeyError: 'netstat'
```

Idx		Hash	Key	Value	
000	=	...11100000	'svn'	3690	x
001	#	...00110100	'dict'	2628	x
010					
011	#	...10111000	'zope'	9673	
100	=	...10001100	'smtp'	21	x
101					
110	#	...01001000	'ircd'	6667	
111					



# Не все поиски одинаковы

- Некоторые находят результат сразу
- Некоторым нужны несколько итераций

```
threes = {3: 1, 3+8: 2, 3+16: 3,
3+24: 4, 3+32: 5}
```

Idx		Hash	Key	Value	
000	≠	...00001011	11	2	
001	≠	...00010011	19	3	✖ ↗
010					
011	=	...00000011	3	1	→ ✖ ↗
100					
101					
110	≠	...00011011	27	4	✖ ↗
111	≠	...00100011	35	5	● ↗

# Удаление элементов

- Нельзя просто так взять, и пометить ячейку как пустую
- Необходимо вставить специальный "dummy" элемент

```
del d['smtp']
```

Idx		Hash	Key	Value
000	=	...11100000	'svn'	3690
001	≠	...00110100	'dict'	2628
010				
011	≠	...10111000	'zope'	9673
100	=	...10001100	'smtp'	21
101				
110	≠	...01001000	'ircd'	6667
111				

```
del d['smtp']
d['ircd'] ???
```

Idx	Hash	Key	Value
→ × 000 = ...11100000	'svn'	3690	
001 ≠ ...00110100	'dict'	2628	×
010			
011 ≠ ...10111000	'zope'	9673	
100 = ...10001100	'smtp'	21	×
101			
110 ≠ ...01001000	'ircd'	6667	●
111			

```
del d['smtp']
```

#Заменяем на "dummy" слот

#Может быть использован снова

Idx		Hash	Key	Value
000	=	...11100000	'svn'	3690
001	≠	...00110100	'dict'	2628
010				
011	≠	...10111000	'zope'	9673
100	!		<dummy>	
101				
110	≠	...01001000	'ircd'	6667
111				

```
del d['smtp']
```

#Заменяем на "dummy" слот

#Может быть использован снова

Idx		Hash	Key	Value	
000	→ × =	...11100000	'svn'	3690	
001	≠	...00110100	'dict'	2628	×
010					
011	≠	...10111000	'zope'	9673	
100	!		<dummy>		×
101					
110	≠	...01001000	'ircd'	6667	●
111					

```
>>> del d['svn'], d['dict'],
d['zope']
>>> d['ircd']
```

	Idx	Hash	Key	Value	
→	x 000 !		<dummy>		✗ ↗
	001 !		<dummy>		✗ ↙
	010				
	011 !		<dummy>		
	100 !		<dummy>		✗ ↙
	101				
	110 ≠ ...01001000	'ircd'	6667		✗ ↙ ↘
	111				

# Увеличение размера таблицы

- Таблица заполнена максимум на 2/3
- 2.7: < 50k size × 4
  - > 50k size × 2
- 3.3: size × 2

```
>>> d = {}
```

Idx	Hash	Key	Value
000			
001			
010			
011			
100			
101			
110			
111			

```
d = dict.fromkeys(words[:5])  
# 40% коллизий  
# Заполнен на 2/3, resize
```

Idx		Hash	Key	Value
000	=	...11100000	'a'	None
001	=	...00100001	'aback'	None
010				
011	=	...00100011	'abaci'	None
100	≠	...00011111	'abase'	None
101				
110				
111	≠	...01110001	'abaft'	None

`d['abash'] = None`

# размер ×4 до 32

# 0% коллизий

00000	= ...011000111100000	'a'	None
00001	= ...110010100100001	'aback'	None
00010			
00011	= ...110010100100011	'abaci'	None
00100			
00101			
00110			
00111			
01000			
01001			
01010			
01011			
01100			
01101			
01110			
01111			
10000			
10001	= ...001101001110001	'abaft'	None
10010	= ...000100100010010	'abash'	None
10011			
10100			
10101			
10110			
10111			
11000			
11001			
11010			
11011			
11100			
11101			
11110			
11111	= ...000100100011111	'abase'	None

```
d = dict.fromkeys(words[:21])
```

# 29% коллизий

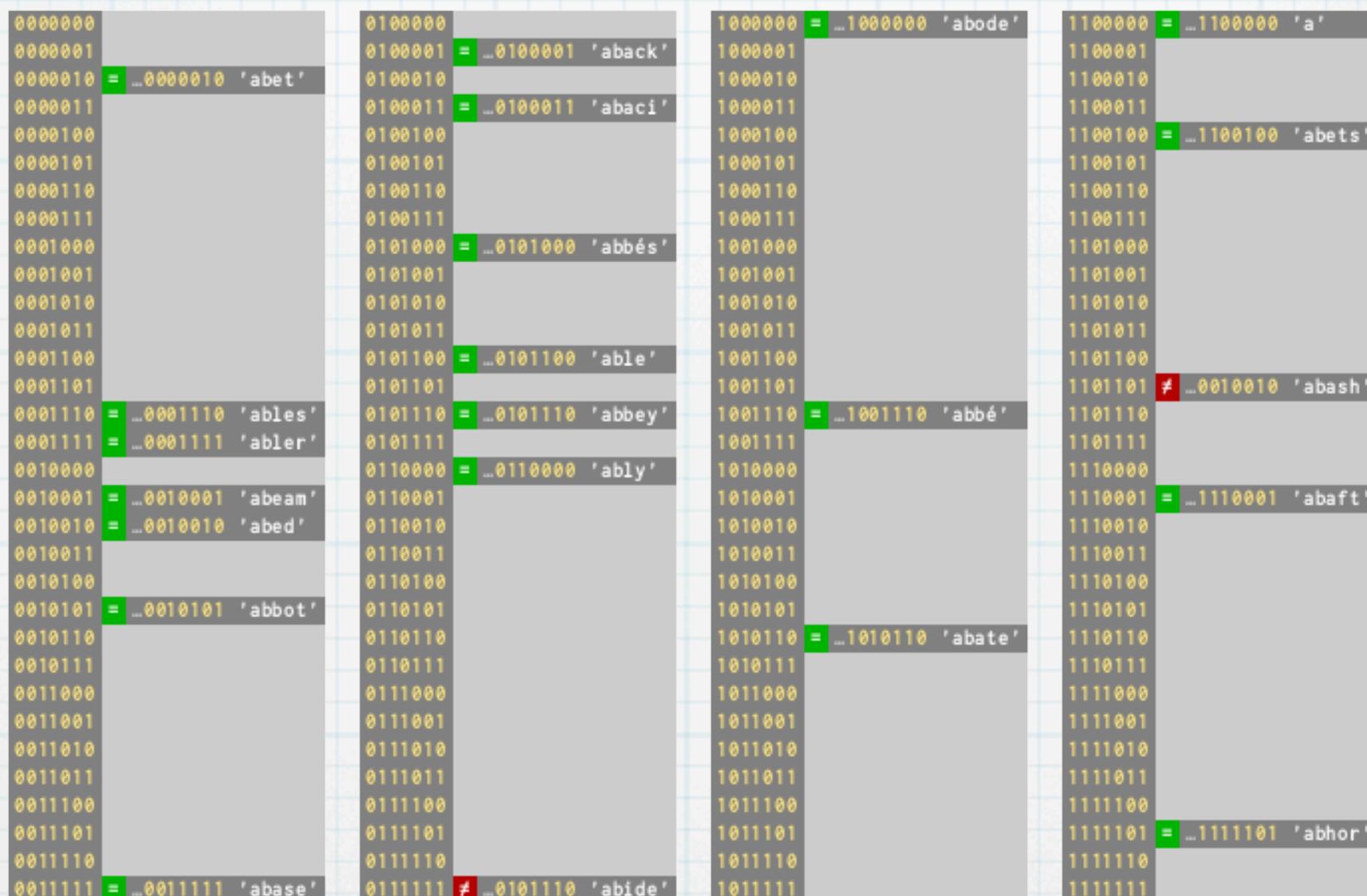
# Заполнен на  $\frac{2}{3}$

00000	=	...011000111100000	'a'	None
00001	=	...110010100100001	'aback'	None
00010	=	...110011100000010	'abet'	None
00011	=	...110010100100011	'abaci'	None
00100	=	...000001011100100	'abets'	None
00101				
00110	#	...011101110110000	'ably'	None
00111	#	...101111110010001	'abeam'	None
01000	=	...001110010101000	'abbés'	None
01001	#	...110110010001110	'ables'	None
01010				
01011				
01100	=	...011101110101100	'able'	None
01101	#	...110011100010010	'abed'	None
01110	=	...111001000101110	'abbey'	None
01111	=	...110110010001111	'abler'	None
10000	#	...001110011001110	'abbé'	None
10001	=	...001101001110001	'abaft'	None
10010	=	...000100100010010	'abash'	None
10011				
10100				
10101	=	...111111110010101	'abbot'	None
10110	=	...100001011010110	'abate'	None
10111				
11000				
11001				
11010				
11011				
11100				
11101	=	...00111111111101	'abhor'	None
11110	#	...000101010101110	'abide'	None
11111	=	...000100100011111	'abase'	None

`d['abode'] = None`

# размер ×4 до 128

# 9% коллизий



```
d = dict.fromkeys(words[:85])
```

# # 33% коллизий

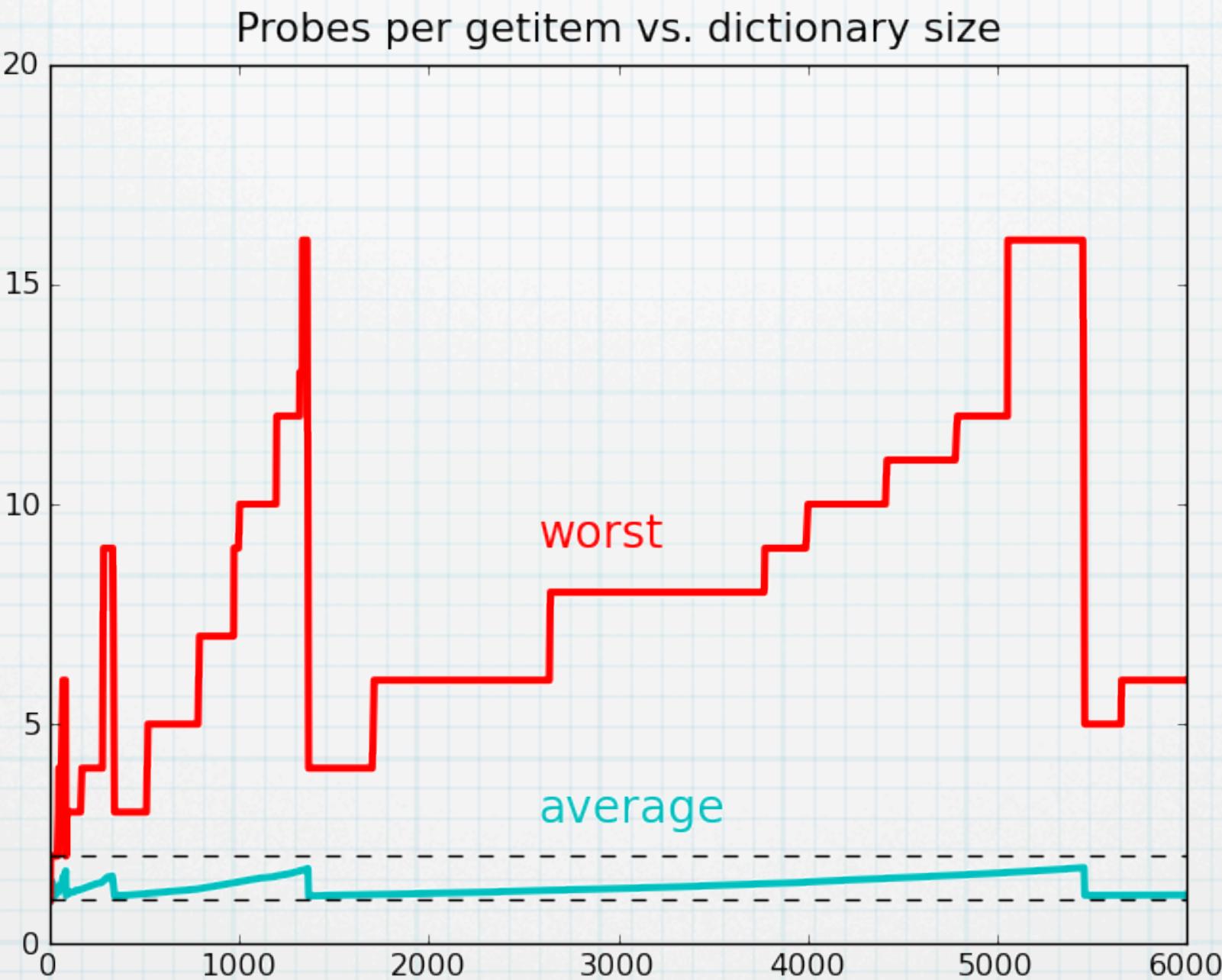
# # Заполнен на 2/3

00000000		01000000	= ...10000000 'abode'	11000000	= ...11000000 'a'
00000001		01000001	# ...010000101 'aback'	11000001	= ...1100001 'acute'
00000010	= ...00000010 'abet'	01000010		11000010	= ...11000010 'about'
00000011	# ...01010111 'acne'	01000011	= ...01000011 'abaci'	11000011	= ...11000011 'adiós'
00000100	= ...00000100 'adobe'	01000100	= ...01000100 'adds'	11000100	= ...11000100 'abets'
00000101	# ...10101110 'acorn'	01000101	# ...10000111 'act's'	11000101	# ...11001110 'acme'
00000110	# ...10101010 'adult'	01000110	# ...01010111 'acts'	11000110	= ...1100110 'ace'
00000111		01000111	# ...01011110 'actor'	11000111	# ...10100001 'adept'
00010000		01010000	= ...01010000 'abbés'	11010000	# ...1000101 'achy'
00010001	# ...01011100 'adz'	01010001		11010001	
00010010	= ...00010100 'acted'	01010010		11010010	
00010011		01010011	# ...01000001 'aches'	11010011	= ...11010111 'acids'
00011000	# ...01000001 'ad's'	01010100	= ...0101100 'able'	11011000	
00011001		01010101	= ...0101101 'ado's'	11011011	# ...00100010 'abash'
00011100	= ...00011110 'ables'	01010110	= ...0101110 'abbey'	11011100	
00011101	= ...00011111 'abler'	01010111	= ...0101111 'achoo'	11011111	# ...11010111 'adopt'
00100000		01100000	= ...01100000 'ably'	11100000	= ...11100000 'adman'
00100001	= ...00100001 'abeam'	01100001	# ...01011111 'acres'	11100001	= ...11100001 'abaft'
00100010	= ...00100010 'abed'	01100010	= ...01100010 'add'	11100010	= ...11100010 'admit'
00100011	= ...00100011 'abort'	01100011		11100011	# ...01011110 'ace's'
00101000	= ...00101000 'aegis'	01101000		11101000	= ...11101000 'abuts'
00101001	= ...00101001 'abbot'	01101001	= ...0110101 'adapt'	11101001	
00101010	# ...00011110 'adzes'	01101010	= ...0110110 'above'	11101010	
00101011		01101011	= ...0110111 'acre'	11101011	= ...1110111 'abuzz'
00110000	= ...00110000 'aces'	01110000		11110000	
00110001		01110001	= ...01110001 'ado'	11110001	# ...1110100 'adore'
00110010		01110010		11110010	
00110011	= ...00110011 'acid'	01110011		11110011	
00111000	= ...00111000 'acrid'	01111000		11111000	
00111001	= ...00111001 'acing'	01111001	= ...0111101 'aeons'	11111001	= ...1111101 'abhor'
00111100		01111010	# ...10000111 'adieu'	11111010	
00111101	# ...00111111 'abase'	01111011	# ...01011110 'abide'	11111011	
		01111100		11111100	
		01111101		11111101	
		01111110		11111110	
		01111111		11111111	

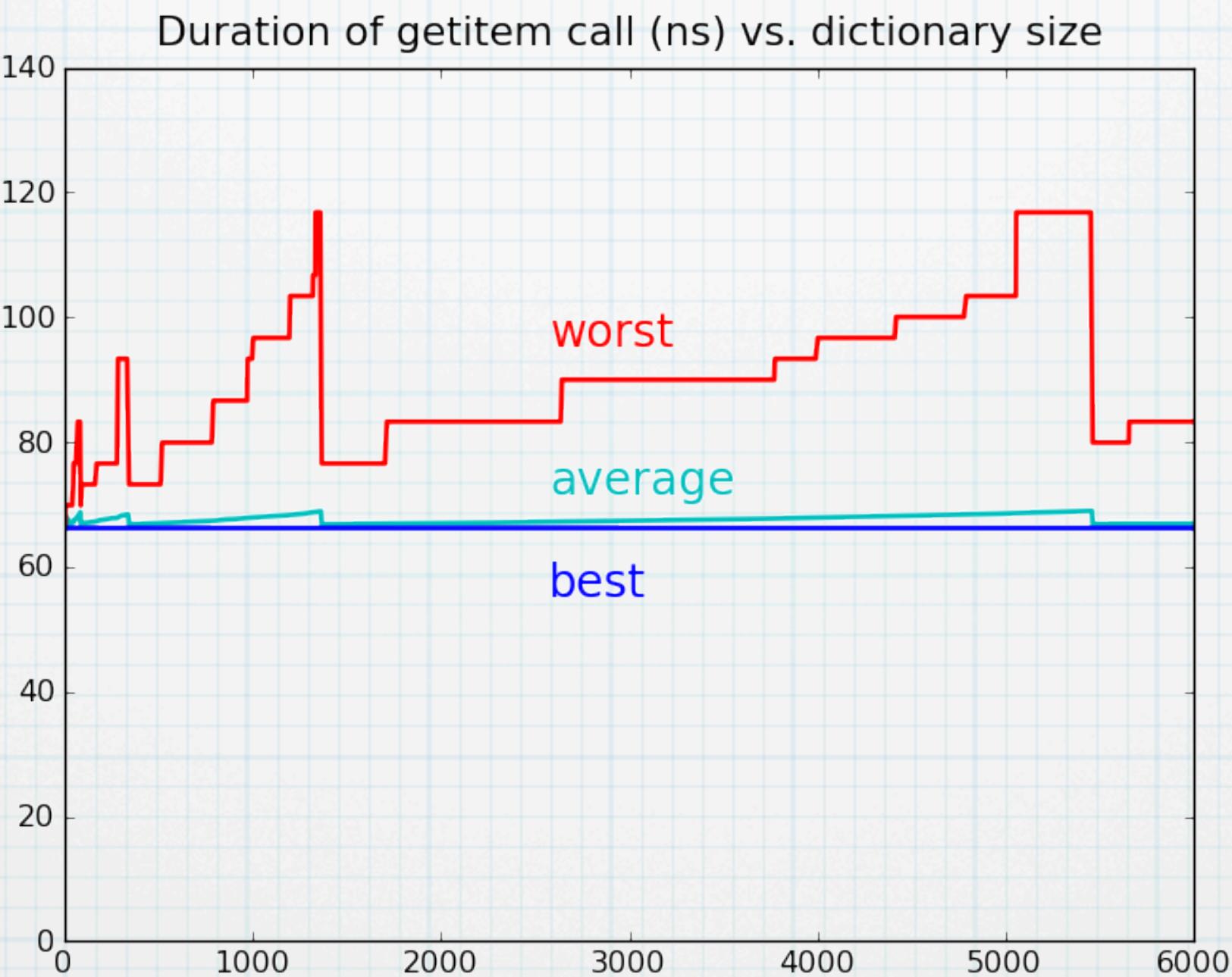
# Время доступа к элементам

- Растет по мере заполнения словаря
- Затем резко уменьшается после изменения размера
- Среднее время доступа O(1)

# Поиски vs размер



# Время vs размер



# Удаление элементов

- Не уменьшает размер таблицы
- Таблица может уменьшиться только при добавлении элементов

# Порядок элементов

- Во время изменения размера порядок элементов может полностью поменяться
- Добавление элементов во время итерации запрещено

**RuntimeError**: dictionary changed size during iteration

# Свой \_\_hash\_\_()

- Хорошо перемешать биты
- Равные хеши для равных элементов
- \_\_eq\_\_() должен быть
- Быстро вычисляется

# Пример \_\_hash\_\_()

```
class Point(object):
    def __init__(self, x, y):
        self.x, self.y = x, y

    def __eq__(self, p):
        return self.x==p.x and self.y==p.y

    def __hash__(self):
        return hash(self.x) ^ hash(self.y)
```

# oCERT #2011-003

- Хэш для str, bytes и datetime смешивается с "солью" уникальной для каждого процесса Python
- pre 3.3: -R option  
3.3: by default

# Python 3: Split-table словари.

## Общая таблица с ключами для разных таблиц со значениями

Idx	Hash	Key	Value	Value
000				
001				
010				
011				
100				
101				
110				
111				

# Спасибо

- <http://blip.tv/pycon-us-videos-2009-2010-2011/pycon-2010-the-mighty-dictionary-55-3352147>
- Python source code