# String

(or, tips and tricks for index design)

# An overview

String

xxxxxxxxxxxxxxxx

xxxxxxxxxxxx

xxginx

# adinadde inelreichddl

## They are ubiquitous:

- iita ibraries a rout ataoues
- etroi hite a yeo aes
- eiaie ioratio soures e eoi or atet bs
-  eb aes reositories
- riate ioratio bs
- 

## tri oetios are roi at a staeri rate:

- ore tha Tb o tetua ata i the eb
- ore tha b o base airs i the eoi bs

# del itled

xntxnick in iimx

xdotxx tin x                                    xx

surfaceu          uauotu u÷uu  u

✓ uu    iions of docents (u    iions per day)u

weepeu          uauotu u u    u

✓ uu u   of interesting text  data

aiing u    istuauotu uu  (eeruar)

✓ uuions of sg per dayuu          ituin u  uaiing ists

# dddrel

dboedthel d           dl

**ouet** is a sie iee o tet otaii soe aru that is seesribioos soe rou rues a is easiy reaabe by huas a outers

✓    t is tet base a ator ieeet        ✓

# decbdriddr d                    d

ueries iht eoit the ta struture to reiera a
seiaie the retriea o the asers or eae:

- **roiity**      ay eoit ta esti

- **or isabiuatio**      ay eoit ta aes



☑ struture is usuay reresete as a          set o aths   stris

☑ queries are ture ito stri queries:          booauthorirstaeaoo

# e need r d d ine d

ruteore sai is ot a iabe aroah:

ast sie searhes

utie sie searhes or oe queries

The eria eritae itioary eies ie as oos

outer siee a ie is a ersistet ata struture that
aos to ous the searh or a query stri or a set o the o
a roaby sa ortio o the ata oetio

# ded     ed d

The ie is a basi bo o ay  syste

syste aso eoasses:

oes

ai aoriths

uery auaes a oeratios

sereeba oes a iteraes

eurity a aess otro aaeet

e i oetrate oy o ie esi

# dd cbcdal clccdel

- ear about:

  oe a raeor or eauati stri ata strutures a aoriths o assie ata sets

    teraeory oe
    auate the oeity o ostrutio a uery oeratios

  ratia a theoretia ouatios o ie esi

    The subsyste a other eory ees
    Tyes o queries a iee ata
    ae s tie traeo
    tri trasatios a ie ahi

  ieeri a eeriets o iteresti iees

    erte ist s ui arrayui tree a tri tree
    o to horeorah oressio a iei: the e rotier

  ✓   : o ear ier ao these ata strutures

deldeadcodeword

# vo id       ie o d n e d d

Tyes o ata

iuisti or toeiabe tet
a sequee o haraters or bytes

sequees
uioieo ies
eutabes

Tyes o query

orbase query
haraterbase query

at or
 or rei or sui
hrase

rbitrary substri
oe athes

To ie i aroahes :

orbase iees          here a oet o or ust be eise

erte iesiature ies or itas

utet iees          o ostrait o tet a queries

ui  rrayui tree ybri ieesor tri tree

String

# chvereldtl  eddr d id

o

o

oabuary

ostis

- **oet o or ust be eise**

  t ees o the ueryi aiatio
  oe squeei: ora orsto orsstei

- **ts sie is usuay sa**

  easa says                $^\beta$here   is the oetio sie
   $\beta$ is ratiay betee  a

- **eetatio**

  rray: ie a sae suit but so queries
  ash tabe: ast eat searhes
  Trie: ast rei searhes but it is ore oiate
   utet ie            ast oe searhes

- **oressio  es           seeu ator o to o sai**

  es ahi a reethi
  eues aout o roesse ata

# del dbon elothd

- **rauarity or        auray i or oatio:**
  - oarseraie        : ee ouet ubers
  - oerateraie            : ee the ubers o the tet bos
  - ieraie        : ee or or setee ubers

ae arou
ast queries a reisio

ae ess tha
o queries: ostiteri

- **orthooa aroah to sae sai:                    a oi**
  - ort the ostis or ireasi ouetbo or ter uber
  - tore the ierees betee aaet osti aues
  - se ariabeeth eois or as:                γoeoob

otiuatio bit: ie bi

✔

padding
agging

t is byteaietaea sesyhroii

ery ast eoi a sa sae oerhea

- Tet oetio iie ito bos o ie sie _____ b

  bo ay sa to or ore ouets
  ostis  bo ubers  ⟵  ingand    ✓ s ✓ againd

- To tyes o sae sais

  utie ourrees i a bo are reresete oy oe
  The uber o bos ay be set to be sa
  ✓ ostis ist is saabout   o the oetio sie
  ✓ er  assae a query tie are o or a roer b

- uery aseri is a threehase roess:

  uery is athe aaist the oabuary:           or athis
  ostis ists o searhe ors are obie:           aiate bos
  aiate bos are eaie to iter out the           ase athes

  oabuary turs    oe    tet searhes ito   eat   bo searhes

  usa or
  suit ie

- **e ostrutio**

  reate oter airs t sorte by ireasi
   eresort o the seo ooet t
  ui ostis ists ro aaet airs ith equa t

  ✓ ae bo eruti or aeotiuous ostis ists

- **ouet uberi**

  oaity i the ostis ists iroes their aoi
  assie eoitatio: teer oi aoriths
   tie eoitatio:              eoreri o o ubers eoh                *et a*

- **atie iei**

  Tas a attributes iee as ters o a roer oabuary
  Ta esti oe as set o              este ri iteras

  ✓ trutura queries ture ito booea a eoetri queries

☑ **ur roet:**  ibraryoressio iei or

**cdllcd** o

- ai iea:

  ereset the ouet tree ia tues or set o obets

  *eetrohere*          ause to aiate ito the tree

  uery eie use staar          *oi*    a    *sa*

  oe aitioa iees or seia aesses

-

## dll d

o

- eera isaataes :

    uery aiatio is ostysiuate ia ay ois

    uery otiiser ooses oee o atura o the ouet

    ies i tabes or shou be sa

    ee etra iees or aai eetie *ath* queries

- rix xi

- rx k Stroonix ritx

# chive drel

The iterature oers arious roosas:

- **set us** : bui a  tree i ai eory at query tie

- **i** : tree or stori airs athor

- **abri** : atriia tree or iei *a* ossibe aths

- **ati** : tree is artitioe ito is aes see e  yee

- **Tey** : tri tree  ✓ are sae ouay

- oe oeria routs : **Taio**  o etais

ae ouay is usuay ot eauate surey it is  ≥

# <span style="color:yellow">idtbbd</span>       <span style="color:cyan">edleelnd</span>

- ouets ay be:

  stroy tetua e iuisti tets

  oy eore a ay our ithout a T

  arbitrariy este a oiate i their ta struture

  retrieabe i their oriia or or  brosers

- 

⟹                                    sasss

# idbbl                    elidd died

- ie ouet iei:

    ie sotare arhiteture

    ustoiabe iei o eah ie they are heteroeeous

    ase o aaeetuate a istributio

    *iht*  itera ie or       *oi ia   tai*           to see u query

- 

-

# eldbod   oel

Their ee is erasie:

a ata        :    sequees uioieo ies

iuisti tets         : ata iistatistis

oabuary      or erte ists

ath      queries o     ouets

trusio etetio tiiruses

## our asses o iees:

✓ ui array or ui tree

✓ Toee iees: ui array  ieory uraie

✓ tree base ata strutures: rei tree

✓ tri tree: tree  atriia trie

ur eture osists o a tour throuh these toos

atter     ours    at ositio i o T

i__ is a rei o the sui Ti

T

Ti

urrees    o i T   suies o T hai as a rei

T Th <u>is</u> is a isua eae

This <u>is</u> a isua eae

This is a <u>isua</u> eae

T    orte set o suies o T

Δ    orte set o suies o a tets i      Δ

ro     suies i  Thai rei  are otiuous

ro     tarti ositio is the eiorahi oe o

Θ    sae

**Sx**    T

T ississii

sui oiter

ui  rray

 : array o its  bytes

Tet T:  bytes

✓    bytes o sae ouay

u  si

irete biary searh o  :                    o            tie

T  ississii

**Sx**

is arer

aesses or biary ste

si

# ithcloodrenedd    cdereболl

ruteore oariso:                          o tie

## T  ississii

**Sx**



si

o

sii                    is a rei

sissii                 is a rei

issii

is ot a rei

ui  rray searh

o              o tie

o              o i ratie

tera eory

ie is ai or

o                        o s

o                        o

o

# adaptive red levd

stores the oestoorei betee suies aaet i

T ississii

xx  Sx  T

ass siss

ui rray searh
o                    o         s
bytes o se

o  {

u  si

oare aaist

S sinsaasas

a uti i

# chde elnedd    d d

reeta searh          usi the  array          : o resai o atter hars

**Sx**

i q

s

s

s

o iut

The ost:   eory aesses

# rddl    ned

oit itera eory         : sae the sui array a oy soethi i eory



u

oy a rei o
are suies

uraie

s o    s    o s

biarysearh isie

Sx                                                is

s

☑ araeter    s ees o         a iuees both erorae a sae

# elidee

cdeid        dl

t is a oate trie          buit o a tet suies

ba        ✓ earh is a ath traersa

tie        a o tie

sae

hat about T i etera eory
baae tree toooy
iaiity                    } ai

are sae        } T tree    o aerae

T abababb

s   s   s   s   s        Ω s
                         Ωo s

sspisnssainsaisand pass

String

# el cbod    oel

e are et ith ay oe issues:

<span style="color:green">ui  rray</span>        : iaiity

<span style="color:green">ui tree</span>        : iiut ai a                    $\Omega$  s

<span style="color:green">ybri:</span>      euristi tui o the erorae


tree is ubiquitous i aresae aiatios:

<span style="color:green">toi eys</span>        : iteersreas

<span style="color:green">rei tree</span>        : boue eth eys            $\leq$   hars


<span style="color:red">ui trees  trees</span>                    <span style="color:red">tri tree</span>          <span style="color:red">erraiarossi</span>

☑ e uboue eth eys

☑ oo orstase  bous i searh a uate

☑ uaratee otia aei ratio

# dedndribnd

tris hae arbitrary eth:

is ae        aot esure the storae o        Θ stris

ay be uabe to store ee oe sie stri

tri storae:

oiters    ao to it      Θ stris er is ae

tri oariso            ees is aess a ay be eesie

tri oiters oraiatio see so ar:

✓ ui array:        sie but stati a ot otia

✓ atriia trie:    sohistiate a uh eiiet otia

ea the robe            : Δ is a tet oetio

✓ earh:                retriee a ourrees o  i            Δs tets

✓ ate Ŧ:                isert or eete a tet Ŧ ro            Δ

# ch dd

tri tree erorae:                                    erraiarossi

earh        taes  o              o s
ateT        taes  t o          s
ae      is Θ  is aes

si the tri tree i itera eory:

earh        taes  o              o tie
ateT        taes  t o          tie
ae      is Θ  bytes
✓ t is a sort o yai sui array

ay other aiatios:

tri sorti                      re     et a
itioary athi                   erraia       et a
utii stri queries              aaish        et a

# Adriidtheerind

dre tri trees aeai i ratie

**String**

gritxxkbrtxx

rotixxtrixx inxignx

xxrginx

# red ichddnidriobnd

ie a tri tree oe                    $\pi$e eie:

$\pi$  set o a stris store at oe                $\pi$

b  aiu sie o

$\pi$

iteresti roerty:

ros as o      b  a oes ot ee o                    $\Delta$ s struture

b is reate to the sae ouay o T            $\pi$a b

☑ The arer is bthe aster are searh a uate oeratios

☑ ur oa: queee        T $\pi$ as uh as ossibe

# dl π id eelnionb

oe π atuay otais et π :

 T π atriia trie iei the stris o

 The oiters to the hire o π

 oe auiiary a booei ioratio


or bytes


eiibe

 the stris are biary the T π ostists o:

 eaesoiti to π s stris

 itera oeseah stori a iteer aue

 arseah stori oe sie har

eeti T π taes:

 bytesia a oiterbase soutio

 bytes ia a roer eoi o the biarytree struture

erraiarossi

# dd did      ddred       d

**eriets hae sho that:**

### earh

t taes about   is aesses            as the orstase bou

t is  ties aster tha ui  rray searh

oarabe to ui  Tree searh

### sertT      ia a bathe isertio

t is  ties aster tha    rei  trees

etter aei ratio tha ui  trees

**To iitatios:**

ae usae o      is too uh

The uate os are  boue

# A new dodd

eeti the oe                    $\pi$:

tri oiters            a    hi oiters         i  bytes

teers        i the oes o T          $_\pi$ store ia otiuatio it

eriets shoe that    are ery sa                    ✓ byte

o o e ieet T                          $_\pi$

hou be sae suit a ao basi aiatioa os

oe resuts o the suit oi o biary trees:

tia        o          bits a basi aiatioa os                    aobso

o          bits a ore aiatioa os                    uro    *et a*

To seiaties o our otet:

T    $_\pi$ is saabout a thousas o stris

aiatioa os  oar traersa

tie is ot the oy resoure                    is surey aie

# T $_\pi$ s toooy ay be roe

Tae the iorer isit o T $_\pi$ :

array o oiters to $_\pi$ s stris ie T $_\pi$ eaes

array o s betee stris aaet i

s    s    s    s    s    s

$_\pi$ s stris o is

# ch dd

oe π otais et π :

oiter array

iteer arraystore by otiuatio it

earhi s ositio ao π s stris:

to eth the is ae otaii oe π

array sas : hars a iteer oarisos

stri aess to the aiate stri s

ie is about a thousas o stris:

The to eth the is ae taes μs

The to array sas are ery ast: μs ahe reethi

The stri aess iht eoy ireeta searh

☑ ae bous as beorea about bytes o sae i ratie !!

# eeldodl          eld

- roie a ubi ieetatio o tri trees
  - ✓ eer to ereey or the

- ath queries:  o to ie a abee tree or ath queries
  - ✓ oauthoraeaoo

- utiiesioa substri queries: utiie reor searh
  - ✓ ay e u eoetri ata strutures i tri trees

- trea o queriesossiby biase:          tri tree is ot otia
  - ✓ ay e eise a seausti ie                          eatorTara

- aheobiious tries:        o eiit araeriatio o
  - ✓ tri tree are baae but eeat

# chelodlndobn

uii a utet ie is a haei tas

# del dnididribnd

e hae areay sho that the ui  rray   a the
orresoi   array suie to bui the tri tree

## o o e bui the arrays   a

- eory aoriths are ieiiet
- ai  t ort eiiet but sae osui    rauser                          *et a*
- ∃ theoretiay otia aorithbut oiate a sae osty

erraia      *et a*

## There eists a aorith hih is                    aeaates      *et a*

- Theoretiay uaetabe: ubi   oeity
- ratiay ery aeai or erorae a sae ouay
- ts asytotis a be iroe ith soe tris                    rauser    *et a*

idArrcdrel  itdd

is

T

T  TTTT  TTTT

eth i eory the irst iee Ta bui  a  or the
   suies hih start at ositios

   ossiby soe etra  s are eee st a  th sui

To is

et

et

utio     :  e hae     et  a     et  or the suies starti isie
   Te ete this to the suies starti i Tii

e ai at eeuti aiy bu  s

# dAArrdrel

chive dd

T

roesse  i

is

et

et

ere

utio____:
eth i eory the et iee Tiia bui a

a To is a oute a ieory outi array

- earh ithi the ositio o eah sui starti ito T

  This taes i s atuay bu s

# iddArrcdlred

ihdiveeld

ere    $_{et}$  a  by usi the array ia a is sa

$_{et}$

$_{et}$

**The oeity o the ith ste is:**

- ethi Tiitaes   s                              bu  s
- uii  a  taes ratiay o  s                      or e raos
- outi  ia a sa o Ttaes i  s                    bu  s
- eri        $_{et}$ia ia taes i  s             bu  s

✓ era the aorith eeutes             s i ratie        aiy bu  s

a sas
isisa issad s

dndbrthd

obrti stris is siiar to sorti suies

d el nde oddinddriahd

itera eorye o a otia bou:

- ia a oate trie e et          Θ o          tie
- oer bou oes ro the sorti o  eeets

etera eorye ou eet to ahiee:

Θ o                    s

but

tri trees ao to ahiee          O o          s

Threeay quisort ets        O o        s              eteyeei

The situatio is uh oiate the oeity ees o                              :

- breai stris ito hars is aoe
- the stri sie reatie to

# elndio

et us eie                                   :

- a        or stris saer tha
- a        or stris oer tha

stris are iiisibe eeryhere              it is otia       :

o                              o

short                      o

stris are oy iiisibe i etera eory:

i        o          o                o

short                    o

stris ay be hoe ito iees:                      s

- t is a raoie aorith                           erraiaThoru
- The aerae stri eth shou be        $\Omega$ o              o

drid

| ab | bb | ab | bb | aa | ab |
|----|----|----|----|----|----|
| ab | b | ab | b | bb | |
| b | a | b | | aa | |
| b | aa | aa | | bb | b |
| ab | bb | bb | aa | aa | ab |

as

s

ssad nass

s nass        ans

aa
ab
bb
a
b

b

sss
issad as

ass

sad an

pssp sas
aangd ssisas

drid drid dnd d

| | | | | | |
|----|----|----|----|----|----|
| ab | bb | ab | bb | aa | ab |
| ab | b | ab | b | bb | |
| b | a | b | | | aa |
| b | aa | aa | | bb | b |
| ab | bb | bb | aa | aa | ab |

ut

ashe a sorte stris

Tabe T ater orar a

ee the surey

s     s

pssp sas
aangd ssias

ssad an

s

s

# eddeld

- ose the arious as
  - ✓ o stris i the ase o iiisibiity o etera eory
  - ✓ etter aaysis or the raoie aorith

- eet a those aoriths

- hat about aheobiious stri sorti aoriths
  - ✓ ost o the are base o tries
  - ✓ rbitrary eth reates a ot o robes
  - ✓ robaby the raoie aroah a he i this ase too

# dcode led d             n eed d

ds sae oerhea the ta to ay or usi a utet ie

# ddlcdeidhcddl

✓ oressio has to ositie eets:

    ✓ ae sai

    ✓ erorae iroeet

        ✓ etter use o eory ees ose to roessor

        ✓ rease is a eory baith

        ✓ eue ehaia see tie

          see aes eoressio ostess

☑   e estabishe    : t is ore eooia to store ata i oresse or tha uoresse

☑  uth i the r o says:      ae otiiatio is osey reate to tie otiiatio i a is eory syste

☑      reease i arh the eory e asio Tehooy Tue ito eerers

    ✓ oube eory at about sae ost a erorae



IBM Journal of Research and Development

# Conclusion

Classia utet iees use $\Theta$ o bits o storae

- ui array: o o tie
- tri tree: o o s

uit sui trees use o $\Theta$ bits o storae

uro et a

eay eee

ui erutatio aot be ay ro

- biary tets erutatios o

oat sui array uses $\Theta$ bits o storae rossiitter

- uery tie is o o o $\varepsilon$ tie

# el dod ed

- ut:

  ostatsie ahabet                    Σ
  arbitrariy o tet Toer                        Σ

- uery   o a arbitrary stri

  out    the ourrees o  i T
  oate    the ositios o the ourrees o  i T

i at eoiti reetitieess i the iut to squeee the ie

oes it eist a oortuisti ie

<u>ae:</u>

- out    the as ro        oe

- oate   ho ae ro et i isa

s

# rii atastruture esi a oressio tehiques:

- ✓ ui array ata struture
- ✓ urros heeer Trasor

✓ bi oressio aorith

# The theoretia resut:

- ✓ uery oeity: o o

$\varepsilon$ tie

- ✓ ae ouay:

T o bits

✓ o i T is oressibe

th orer eiria etroyit ay be o

# The ie stu is that this resut:

- ✓ is ieeet o the iut soureie oitise o T
- ✓ iiitey shos that ui rrays are oressibe

# ratiethe ie is uh aeai:

- ✓ ae ose to the best o oressors
- ✓ uery tie o e iises o hures o s o tet

# elodrd

et us ie a tet T ississii

ort the res

T

ery ou is a erutatio o Thee aso a

# dUnverdl e dnd d

To roerties:

- s hars reee        s i T
- ith i     ith i

eostrut T
i   tie

T                    i

# dldd dlde idl e

To obseratios:

- qua substr rei aaet ros
- ose hars are siiar

oaity

orith i :

- ✓ oetorot oi o          ✓
- ✓ ueth oi o              ✓
- ✓ tatistia oer o          : ritheti

☑ i oresses uh better tha ibut it soer i eoressio

d  cddn d

si

i

irst ste

s

e

o
pps
s
e

utie ste:        ie   s e ori

✓  Tae i

✓  i irst i s

✓  i ast i s

✓  to ai o these hars

# dd d ocrendd

si

T ississii

s   s   s   s

sai ste is

theorye set to o
baae sae a isti tie

ε to

or thise ee to o baar:

ississis sii

s

e s ro

ro    s s ositio e et                o

s ⇒

e ⇒

This ourree is iste ieiatey

# eldlelin dcl

e eeoe to toos:

- ✓ Tiy ie    suorts ust the outi o the ourrees
- ✓ at ie    suorts both out a oate

both o the easuate a oresse oy o the tet

| uoaction u u neu s uuu | spas | ssiss |
|---|---|---|
| ins  ie | | s<br>outi oy |
| as  ie | | s<br>oati oe ourree |
| sp  o sp e ies<br>                 sps | | s |

☑ *uossass fingerprint* u uxistentiaund conting u  eries fast

# daddblddnedd

hat about orbase ourree o

- ✓ earh or as a substri o Tusi the ie
- ✓ or eery aiate ourreehe i it a orbase oe

or     rei     substri     sui        bi

T bi bi ubi ubi ◯    ◯    ◯

the ostroessi hase a be ery osty

The ie a be aate to be orbase

- ✓ reroess T to or a ieste tet T
- ✓ ui a ie oer T
- ✓ Trasor ay orbase ourree o Tito a substri
  ourree o Ta soe it usi the ie buit o T

# elided

ariat o ua aorith:

- ✓ ybos o the hua tree are the ors o T
- ✓ The ua tree has aout
- ✓ oeors are byteaie a tae

y or

| α | ω | ρ |

| α | ω | ρ |

bi

| α | β |

T bi or ot bi

yes
o

T

| α | β | β | γ | α | β |

| β | γ | γ | α | β | α | β |

yes

o

ie
itioary o ors

ua tree

ie buit o T

passss
sd asssss

# eddded

- hiee o tie i ourree retriea
  - ✓ T o ε o bits          erraia aii

- hiee o s i ourree retriea
  - ✓ o oresse ie es eror rao aesses

- ast ostutio aoriths or ui rrays
  - ✓ i oressio or ie ostrutio
  - ✓ ui Tree ostrutio
  - ✓ usteri o ouets

- eet the too:   ie ise
  - ✓ This iroes theoretiay the erte ists

# d end

y e yearse i be abe to        store eerythi

                                        ray

ato   i   *haerus*   sueste that   riti   ou rate

oretuess i the is o those ho ear to use it
a the sho o iso ithout the reaity

hoe that this i ot ouraai