

Defⁿ: A GNFA is a 5-tuple $N = (Q, \Sigma, \delta, s, F)$, where now δ is a set of generalized instructions (p, α, q) where α is a regexp over Σ .

Note: An NFA_ϵ is the same as a GNFA in which α is a basic regexp ϵ or Σ (or \emptyset for non-edges).

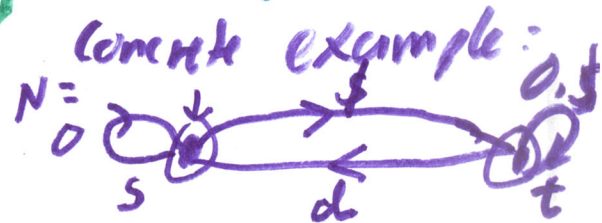
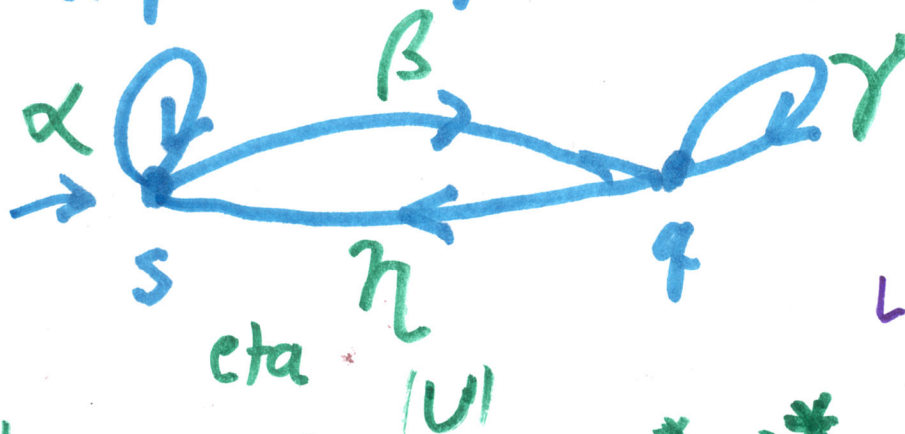
Defⁿ: A valid computation trace of N on X has the form $(p, u_1, q_1, u_2, q_2, u_3, \dots, q_{m-1}, u_m, q)$

where $u_1 \dots u_m = X$ and for each j , $1 \leq j \leq m$, there is an instruction (q_{j-1}, α_j, q_j) st. $u_j \in L(\alpha_j)$.

Then we say $X \in \underline{L_{pq}}$, and $L(N) \stackrel{\text{def}}{=} \bigcup_{f \in F} L_{sf}$ as before.

Note: For every regexp α we can build the trivial equivalent GNFA $s \xrightarrow{\alpha} f$. The conversion from GNFA to regexp has this as the end case.

Examples and Practical Time-Saver: 2-state NFA



$$L(N) = L_{ss} \cup L_{st}$$

$$= L_{ss} \cdot (\epsilon + \$ (0 + \$)^*)$$

where

$$L_{ss} = (0 + \$ (0 + \$)^* d)^*$$

This explains the regular expression in the "Turning Kit" "DragonSh" example

$$L_{ss} = (\alpha + \beta \gamma^* \eta)^*$$

$$L_{sq} = L_{ss} \cdot \beta \gamma^*$$

$$\text{also} = \alpha^* \beta \cdot L_{qq} = \alpha^* \beta \cdot (\gamma + \eta \alpha^* \beta)^*$$

Without adding an extra start and final state, we can use this as a base case whenever the original N has at most 1 acc. state besides s different from s

Then, let us assume states 3, ..., k are nonaccepting

Algm: For $k = k$ down to 3:

Bypass and eliminate state k by:

for each incoming edge (i, β, k) :

for each outgoing edge (k, η, j) :

update any existing (i, α, j)

to $\alpha + \beta \gamma^* \eta$ where (k, γ, k) is the loop at k

Extra: In UNIX, $\alpha? \equiv \alpha + \epsilon$. Best used only with single chars and other short exprs, eg $c?$ or $[a-z]?$ for an optional lowercase letter.

IMHO the precedence of $?$ is hard to remember - whereas the precedence of $+$, \circ , $*$ is the same as for plus, times, and exponents in math.

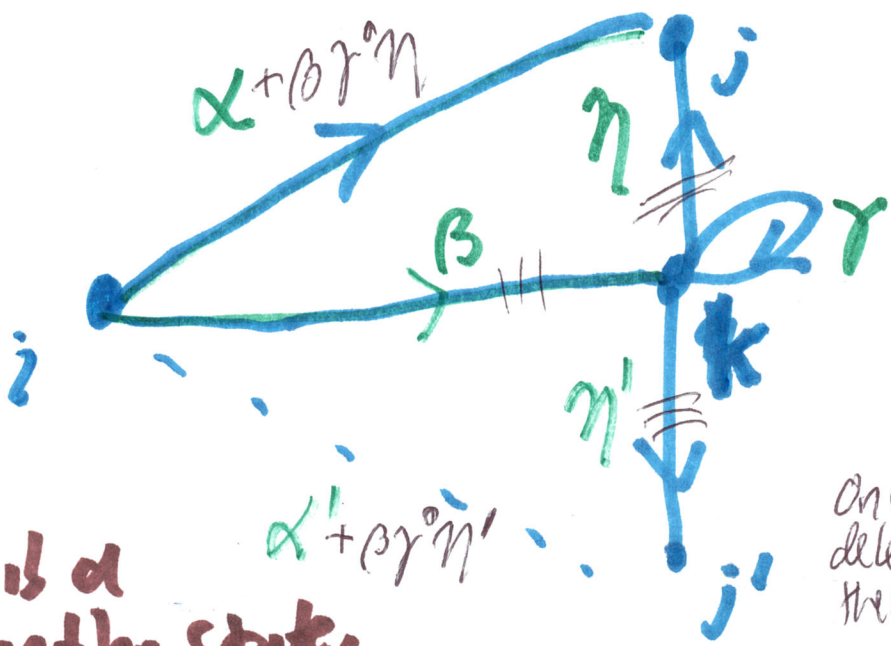
Note how this lecture showcased nested $*$ s as in $(ba^*a + aa^0(b+aa^0))$. The combination of allowing \sim as an extra operator on nested $*$ parts is what can blow up the running time of egrep: extended regular expressions. But if negation is limited to char level (as in $[^a-z]$ for not a lowercase char) then it doesn't blow up. Several NP-hard problems become tractable when a parameter related to allowable nesting is held to a fixed bound.

"Fixed-parameter tractable"

This also shows why although the regular languages are closed under \sim , we hesitate to allow it as a basic operation. Our brains are wired for \cup , \circ , and $*$, but not for \sim (and sometimes not for \cap either).

Picture

$\alpha = \emptyset$ possible



$\gamma = \emptyset$ possible then $\gamma^* = \epsilon$.
 Since $\emptyset^* = \epsilon$.
 $\beta \neq \emptyset, \eta \neq \emptyset$.

once all outedges are handled delete all incoming edges and then delete k .

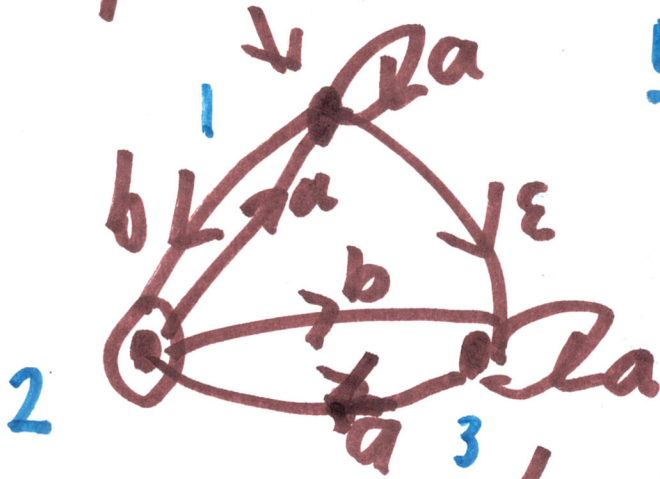
Since k is a non-accepting state, any processing that goes in to k from some node i must come out at some node j , having matched $\beta\gamma^*$. It could equivalently match $\alpha + \beta\gamma^n$ directly.

In practice, can shortcut examples (4)

Elim 3:

Inc: (1, 2, 3), (2, b, 3)

Out: (3, a, 2)



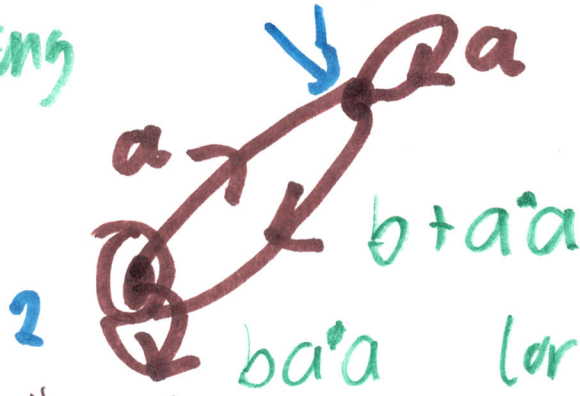
\therefore update (1, -, 2), (2, -, 2)

$$(1-2)_{\text{new}} = (1-2)_{\text{old}} + (1-3)(3-3)^* (3-2)$$

$$= b \quad \epsilon \quad a^* \quad a$$

$$(2-2)_{\text{new}} = (2-2)_{\text{old}} + (2-3)(3-3)^* (3-2)$$

or $\epsilon \rightarrow \emptyset$ + b a^{*} a
won't be wrong



Final answer: $L_{12} = a^*(b + a^*a) \cdot L_{22}$ (or = $\epsilon + ba^*a$, huh?)

where $L_{22} = (ba^*a + a a^*(b + a^*a))^*$

An extra ϵ inside the outer * would not change L_{22} and the final answer.