

Wildcard Matching in the Spirit of E. W. Dijkstra

Herwig Egghart, Scientific Games International

June 25, 2015

Let a, b, \dots denote symbols from some alphabet Σ , and let $*$ be a special symbol outside Σ .
 s, t, \dots shall denote finite strings over Σ , and p, q, \dots shall denote finite strings over $\Sigma \cup \{*\}$.

Then we define the wildcard-matching relation \sim between "sample" and "pattern" strings by the following little calculus:

$$(0) \vdash s \sim s$$

$$(1) \vdash s \sim *$$

$$(2) s \sim p, t \sim q \vdash s \cdot t \sim p \cdot q$$

The two axioms (0) and (1) capture our intuition that any string matches itself, and any string matches a star. Rule (2) expresses the compatibility between wildcard matching and string concatenation (indicated by a dot-shaped infix operator, which is non-symmetric and associative).

Is this calculus a good basis to derive a program that decides whether $s \sim p$?

No, because it only enumerates successful matches, while our program will also have to detect unsuccessful matches. Therefore we shall rephrase the properties of \sim in a more systematic way based on all possible strings on both sides of the relation.

The sample (LHS) can either be the empty string ϵ , or some arbitrary symbol a followed by a tail string s . The pattern (RHS) can either be ϵ , or some non-star symbol b followed by tail p , or a star symbol followed by tail p . If we combine the two LHS and the three RHS cases, we obtain the following six laws:

$$(3) \epsilon \sim \epsilon \equiv \text{true}$$

$$(4) a \cdot s \sim \epsilon \equiv \text{false}$$

$$(5) \epsilon \sim b \cdot p \equiv \text{false}$$

$$(6) a \cdot s \sim b \cdot p \equiv a = b \wedge s \sim p$$

$$(7) \epsilon \sim * \cdot p \equiv \epsilon \sim p$$

$$(8) a \cdot s \sim * \cdot p \equiv a \cdot s \sim p \vee s \sim * \cdot p$$

Although most of these laws are fairly obvious, the total ensemble (3 - 8) is not as simple as calculus (0 - 2) any more. Therefore we had better convince ourselves that both definitions are equivalent, i.e. that the calculus is sound and complete with respect to the laws.

Let us start with soundness, and more specifically with the two axioms. Can (0) and (1) be derived from (3 - 8) for arbitrary s ? Well, s is either empty or non-empty. (0) for empty s evaluates to:

$$\begin{aligned} & \epsilon \sim \epsilon \\ = & \{ (3) \} \text{ Remark } \text{ The curly brackets contain a hint to justify the proof step. (End of Remark.)} \\ & \text{true} \end{aligned}$$

Now let s be non-empty, composed of a and t , say:

$$\begin{aligned}
 & a \cdot t \sim a \cdot t \\
 = & \{ (6) \} \\
 & a = a \wedge t \sim t \\
 = & \{ \text{drop true conjunct} \} \\
 & t \sim t
 \end{aligned}$$

This looks much like our original proof obligation $s \sim s$, except for the different letter! However, since t is always one symbol shorter than s , repeated application of this proof step eventually leads to $s = \epsilon$, for which (0) holds as we have seen. Thus axiom (0) is sound.

Similarly, we observe for axiom (1) and empty s :

$$\begin{aligned}
 & \epsilon \sim * \\
 = & \{ (7) \text{ with } p := \epsilon \} \\
 & \epsilon \sim \epsilon \\
 = & \{ (3) \} \\
 & \text{true}
 \end{aligned}$$

Otherwise let $s = a \cdot t$ again:

$$\begin{aligned}
 & a \cdot t \sim * \\
 = & \{ (8) \text{ with } p := \epsilon \} \\
 & a \cdot t \sim \epsilon \vee t \sim * \\
 = & \{ (4) \} \\
 & \text{false} \vee t \sim * \\
 = & \{ \text{drop false disjunct} \} \\
 & t \sim *
 \end{aligned}$$

Like with axiom (0), we have gotten rid of one symbol in this proof step. So due to the same induction argument, axiom (1) is sound as well.

To prove the soundness of rule (2), we have to show

$$(9) \quad s \sim p \wedge t \sim q \Rightarrow s \cdot t \sim p \cdot q \quad (\text{concatenation lemma})$$

for arbitrary samples and patterns.

The consequent $s \cdot t \sim p \cdot q$ calls for a case distinction based on the structure of s and p , because they appear at the beginning of $s \cdot t$ and $p \cdot q$, and our laws work from left to right. So we have the six cases we already know from the construction of the six laws.

Case 3:

$$\begin{aligned}
 & \epsilon \sim \epsilon \wedge t \sim q \\
 = & \{ (3) \} \\
 & \text{true} \wedge t \sim q \\
 = & \{ \text{drop true conjunct} \} \\
 & t \sim q \\
 = & \{ \epsilon \text{ unit element} \} \\
 & \epsilon \cdot t \sim \epsilon \cdot q
 \end{aligned}$$

Case 4:

$$\begin{aligned} & a \cdot s \sim \varepsilon \wedge t \sim q \\ = & \{ (4) \} \\ & \text{false} \wedge t \sim q \\ \Rightarrow & \{ \text{ex falso quodlibet} \} \\ & a \cdot s \cdot t \sim \varepsilon \cdot q \end{aligned}$$

Case 5:

$$\begin{aligned} & \varepsilon \sim b \cdot p \wedge t \sim q \\ = & \{ (5) \} \\ & \text{false} \wedge t \sim q \\ \Rightarrow & \{ \text{ex falso quodlibet} \} \\ & \varepsilon \cdot t \sim b \cdot p \cdot q \end{aligned}$$

Case 6:

$$\begin{aligned} & a \cdot s \sim b \cdot p \wedge t \sim q \\ = & \{ (6) \} \\ & a = b \wedge s \sim p \wedge t \sim q \\ \Rightarrow & \{ \text{induction hypothesis} \} \\ & a = b \wedge s \cdot t \sim p \cdot q \\ = & \{ (6) \} \\ & a \cdot s \cdot t \sim b \cdot p \cdot q \end{aligned}$$

Case 7:

$$\begin{aligned} & \varepsilon \sim * \cdot p \wedge t \sim q \\ = & \{ (7) \} \\ & \varepsilon \sim p \wedge t \sim q \\ \Rightarrow & \{ \text{induction hypothesis} \} \\ & \varepsilon \cdot t \sim p \cdot q \\ \Rightarrow & \{ (7) \text{ or } (8) \text{ depending on } t \} \\ & \varepsilon \cdot t \sim * \cdot p \cdot q \end{aligned}$$

Case 8:

$$\begin{aligned} & a \cdot s \sim * \cdot p \wedge t \sim q \\ = & \{ (8) \} \\ & (a \cdot s \sim p \vee s \sim * \cdot p) \wedge t \sim q \\ = & \{ \wedge \text{ distributes over } \vee \} \\ & (a \cdot s \sim p \wedge t \sim q) \vee (s \sim * \cdot p \wedge t \sim q) \\ \Rightarrow & \{ \text{induction hypothesis twice} \} \\ & a \cdot s \cdot t \sim p \cdot q \vee s \cdot t \sim * \cdot p \cdot q \\ = & \{ (8) \text{ with } s, p := s \cdot t, p \cdot q \} \\ & a \cdot s \cdot t \sim * \cdot p \cdot q \end{aligned}$$

Hence the concatenation lemma (9) always holds, and so the whole calculus (0 - 2) is sound.

Now to the completeness part: Can the calculus produce every wildcard match that holds according to laws (3 - 8)? For law (3), the answer is easy:

$$\begin{array}{l} \vdash \{ (0) \text{ with } s := \varepsilon \} \\ \varepsilon \sim \varepsilon \end{array}$$

Laws (4) and (5) only specify false matches - which cannot be produced as the calculus is sound - so there is nothing left to show.

Laws (6 - 8) specify conditional matches, with all RHS occurrences of \sim ($s \sim p$; $\varepsilon \sim p$; $a \cdot s \sim p$, $s \sim * \cdot p$) containing fewer symbols than the corresponding LHS occurrences ($a \cdot s \sim b \cdot p$; $\varepsilon \sim * \cdot p$; $a \cdot s \sim * \cdot p$). Hence we can use \vdash RHS as induction hypothesis and only need to show that \vdash LHS follows, ideally by showing RHS \vdash LHS.

For law (6):

$$\begin{array}{l} a = b, s \sim p \\ \vdash \{ (0) \text{ with } s := a \} \\ a \sim b, s \sim p \\ \vdash \{ (2) \} \\ a \cdot s \sim b \cdot p \end{array}$$

For law (7):

$$\begin{array}{l} \varepsilon \sim p \\ \vdash \{ (1) \text{ with } s := \varepsilon \} \\ \varepsilon \sim *, \varepsilon \sim p \\ \vdash \{ (2) \} \\ \varepsilon \sim * \cdot p \end{array}$$

Law (8) requires two proofs, because any of the conditions $a \cdot s \sim p$, $s \sim * \cdot p$ is sufficient. Let us start with the easier part by showing $a \cdot s \sim p \vdash a \cdot s \sim * \cdot p$:

$$\begin{array}{l} a \cdot s \sim p \\ \vdash \{ (1) \text{ with } s := \varepsilon \} \\ \varepsilon \sim *, a \cdot s \sim p \\ \vdash \{ (2) \} \\ a \cdot s \sim * \cdot p \end{array}$$

Can we show $s \sim * \cdot p \vdash a \cdot s \sim * \cdot p$ as well? Unfortunately not, because our calculus has no mechanism for prepending a symbol to the sample alone. If p is empty, however, we have:

$$\begin{array}{l} s \sim * \\ \vdash \{ (1) \text{ with } s := a \cdot s \} \\ a \cdot s \sim * \end{array}$$

And if p is non-empty, the only possible source of $\vdash s \sim * \cdot p$ is an application of (2) with $* \cdot p$ emerging from some shorter prefix $* \cdot q$ and a non-empty suffix r . At the same time, s emerges from parts \dagger and u , say:

$$\dagger \sim * \cdot q, u \sim r \vdash s \sim * \cdot p$$

Now suppose $(\vdash t \sim *q) \Rightarrow (\vdash a \cdot t \sim *q)$ for all q shorter than p .

Then we can prove $(\vdash s \sim *p) \Rightarrow (\vdash a \cdot s \sim *p)$ by the following induction step:

$$\begin{aligned}
 & \vdash s \sim *p \\
 \Rightarrow & \{ \text{undo application of (2)} \} \\
 & \vdash t \sim *q, u \sim r \\
 \Rightarrow & \{ \text{induction hypothesis} \} \\
 & \vdash a \cdot t \sim *q, u \sim r \\
 \Rightarrow & \{ (2) \} \\
 & \vdash a \cdot t \cdot u \sim *q \cdot r \\
 = & \{ s = t \cdot u \wedge *p = *q \cdot r \} \\
 & \vdash a \cdot s \sim *p
 \end{aligned}$$

Together with the induction basis $(\vdash s \sim *\epsilon) \Rightarrow (\vdash a \cdot s \sim *\epsilon)$ already proved before, this discharges our last proof obligation for the completeness of calculus (0 - 2). So not only do the six laws agree with everything the calculus produces, but the calculus is also able to produce whatever the laws claim to be true. Thus, laws (3 - 8) are an appropriate basis for our intended programming task.

* * *

Let s and p be the input strings we try to match, and let y and n be two Boolean outputs indicating whether $s \sim p$ ("yes" or "no"). The idea is to initialize both flags with false and wait till the right answer is known. Then we shall change one of the two flags to true.

To make sure that y and n refer to the original values of s and p , we introduce the rigid variable M for the matching predicate $s \sim p$ before program execution. With these ingredients, our program can be specified as follows:

$$\{ M \Xi s \sim p \} \text{ PROG } \{ (y \Xi M) \wedge (n \Xi \neg M) \}$$

So y and n depend on M , which is of course invisible to the program.

But maybe we can keep $M \Xi s \sim p$ invariant while modifying s and p , thus allowing the program to set the correct flag nevertheless. The invariant should also include a weakened form of the postcondition, which is easily achieved by replacing Ξ with \Rightarrow and initializing both flags with false.

$$\begin{aligned}
 & \text{PROG:} \\
 & \{ M \Xi s \sim p \} \\
 & y, n := \text{false}, \text{false;} \\
 \text{I: } & \{ (M \Xi s \sim p) \wedge (y \Rightarrow M) \wedge (n \Rightarrow \neg M) \} \\
 & \text{do } y = n \rightarrow \{ I \wedge \neg y \wedge \neg n \} \\
 & \quad \text{BODY } \{ I \} \\
 & \text{od} \\
 & \{ I \wedge y \neq n \} \\
 & \{ (y \Xi M) \wedge (n \Xi \neg M) \}
 \end{aligned}$$

Now is the time to exploit the six laws of wildcard matching! Depending on the structure of s and p , the loop body should either set a flag or shorten at least one string under invariance of I . Either way, the number of false flags plus the total string length will decrease, which ensures termination of the loop. (We have already seen that strings get shorter from LHS to RHS in the completeness proof.)

In the following text, the functions **H** and **T** ("head" and "tail") shall decompose a non-empty string into its first symbol and the rest.

BODY:

```

{ I ∧ ¬y ∧ ¬n }
if p = ε →
    if s = ε → FRAG3
    [] s ≠ ε → FRAG4
fi
[] p ≠ ε →
    if H.p ≠ * →
        if s = ε → FRAG5
        [] s ≠ ε → FRAG6
        fi
    [] H.p = * →
        if s = ε → FRAG7
        [] s ≠ ε → FRAG8
        fi
    fi
fi
{ I }

```

Note that the six fragments FRAG3 to FRAG8 correspond to the laws (3 - 8).

FRAG3:

```

{ I ∧ s = ε ∧ p = ε }
{ I ∧ M }
y := true

```

FRAG4:

```

{ I ∧ s ≠ ε ∧ p = ε }
{ I ∧ ¬M }
n := true

```

FRAG5:

```

{ I ∧ s = ε ∧ H.p ≠ * }
{ I ∧ ¬M }
n := true

```

FRAG6:

```

{ I ∧ s ≠ ε ∧ H.p ≠ * }
{ I ∧ (M ≡ H.s = H.p ∧ T.s ~ T.p) }
if H.s = H.p →
    { I ∧ (M ≡ T.s ~ T.p) }
    s, p := T.s, T.p
[] H.s ≠ H.p →
    { I ∧ ¬M }
    n := true
fi

```

FRAG7:

$$\{ I \wedge s = \varepsilon \wedge H.p = * \}$$

$$\{ I \wedge (M \Xi s \sim T.p) \}$$

$$p := T.p$$

So far, everything went smoothly because \sim appears at most once on the RHS of laws (3 - 7). By contrast, the RHS of law (8) is a disjunction of two \sim terms, asserting $M \Xi s \sim T.p \vee T.s \sim p$ at the beginning of FRAG8.

But which of these alternatives should we pursue first? Skip the star at $H.p$ and try to match s against $T.p$? Or skip $H.s$ and still face a star to be matched? Well, the latter choice only perpetuates the problem; and its repeated application would discard the whole sample without learning anything new.

So let us first look beyond the star and compare the beginning of $T.p$ with the beginning of s instead. Note that to avoid a premature commitment, we shall leave s and p unchanged and just look for a common prefix of s and $T.p$:

$$s = t \cdot u \wedge p = * \cdot t \cdot q$$

Of course, this is trivial to establish for the empty prefix ($t = \varepsilon$), but maybe we can find the greatest common prefix by "pushing" t deeper into u and q , and then know enough to re-establish I .

FRAG8:

$$\{ I \wedge s \neq \varepsilon \wedge H.p = * \}$$

$$t, u, q := \varepsilon, s, T.p;$$

$$J: \{ I \wedge s \neq \varepsilon \wedge s = t \cdot u \wedge p = * \cdot t \cdot q \}$$

$$\text{PUSH;}$$

$$\{ J \wedge t \text{ maximal} \}$$

$$\text{FINISH}$$

$$\{ I \}$$

For the design of PUSH, we calculate the weakest precondition of the smallest possible push step and invariant J:

$$\text{wp. } (t, u, q := t \cdot H.u, T.u, T.q) . J$$

$$= \{ \text{axiom of assignment} \}$$

$$I \wedge s \neq \varepsilon \wedge s = t \cdot H.u \cdot T.u \wedge p = * \cdot t \cdot H.u \cdot T.q$$

$$= \{ \text{unify with J} \}$$

$$J \wedge u = H.u \cdot T.u \wedge q = H.u \cdot T.q$$

$$= \{ \text{definition of H and T} \}$$

$$J \wedge u \neq \varepsilon \wedge q \neq \varepsilon \wedge H.u = H.q$$

Not surprisingly, we can keep pushing as long as we haven't run out of symbols and the next symbol agrees in both strings. Before we cast our findings into a loop, however, we remember that undefined sub-expressions are forbidden in a guard and therefore introduce a Boolean flag d to indicate that a difference has been found.

```

PUSH:
  { J }
  d := false;
K: { J ∧ (d ⇒ H.u ≠ H.q) }
  do u ≠ ε ∧ q ≠ ε ∧ ¬d →
    if H.u = H.q →
      { J ∧ ¬d ∧ u ≠ ε ∧ q ≠ ε ∧ H.u = H.q }
      t, u, q := t·H.u, T.u, T.q
      { J ∧ ¬d }
    [] H.u ≠ H.q →
      { J ∧ (true ⇒ H.u ≠ H.q) }
      d := true
      { K }
    fi
  { K }
od
{ K ∧ (u = ε ∨ q = ε ∨ d) }
{ J ∧ (u = ε ∨ q = ε ∨ H.u ≠ H.q) }

```

As a by-product of the last derivation, we have also obtained a precise criterion for the maximality of t , viz. that u or q are empty or start with a different symbol.

Now that we have found the common prefix t , let us see if we can generalize law (6) from the elimination of a common head symbol to the elimination of a common prefix:

$$(10) \ t \cdot u \sim t \cdot q \equiv u \sim q \quad (\text{prefix lemma})$$

For empty t , LHS and RHS are equal already; and if we prepend the symbol a to t , we obtain the following induction step as a proof of (10):

$$\begin{aligned}
& a \cdot t \cdot u \sim a \cdot t \cdot q \\
= & \{ (6) \} \\
& t \cdot u \sim t \cdot q \\
= & \{ \text{induction hypothesis} \} \\
& u \sim q
\end{aligned}$$

Back to the postcondition of PUSH, which is the precondition of the missing part FINISH. By far the most interesting case is when we PUSH into another star, i.e. q has the shape $* \cdot r$. So how can we compute $t \cdot u \sim * \cdot t \cdot * \cdot r$?

Obviously $u \sim * \cdot r$ is a sufficient condition for $t \cdot u \sim * \cdot t \cdot * \cdot r$; all we need is calculus (0 - 2):

$$\begin{aligned}
& u \sim * \cdot r \\
\vdash & \{ (0) \} \\
& t \sim t, u \sim * \cdot r \\
\vdash & \{ (1) \} \\
& \varepsilon \sim *, t \sim t, u \sim * \cdot r \\
\vdash & \{ (2) \text{ twice} \} \\
& t \cdot u \sim * \cdot t \cdot * \cdot r
\end{aligned}$$

But is it also a necessary condition? Again, we start with empty \dagger and try to prove:

$$u \sim **r \Rightarrow u \sim *r$$

This is trivial for $u := \varepsilon$:

$$\begin{aligned} & \varepsilon \sim **r \\ = & \{ (7) \text{ with } p := *r \} \\ & \varepsilon \sim *r \end{aligned}$$

Now the induction step $u := a \cdot u$:

$$\begin{aligned} & a \cdot u \sim **r \\ = & \{ (8) \text{ with } s, p := u, *r \} \\ & a \cdot u \sim *r \vee u \sim **r \\ = & \{ \text{induction hypothesis} \} \\ & a \cdot u \sim *r \vee u \sim *r \\ = & \{ (8) \text{ with } s, p := u, r \} \\ & a \cdot u \sim r \vee u \sim *r \vee u \sim *r \\ = & \{ \vee \text{ idempotent} \} \\ & a \cdot u \sim r \vee u \sim *r \\ = & \{ (8) \text{ backward} \} \\ & a \cdot u \sim *r \end{aligned}$$

Note that the last two calculations consist only of equivalence transformations, so we have even proved the stronger lemma:

$$(11) \quad u \sim **r \equiv u \sim *r \quad (\text{eclipse lemma})$$

With this as induction basis, we shall now prepend another symbol a to the common prefix \dagger :

$$\begin{aligned} & a \cdot \dagger \cdot u \sim * \cdot a \cdot \dagger \cdot **r \\ = & \{ (8) \text{ with } s, p := \dagger \cdot u, a \cdot \dagger \cdot **r \} \\ & a \cdot \dagger \cdot u \sim a \cdot \dagger \cdot **r \vee \dagger \cdot u \sim * \cdot a \cdot \dagger \cdot **r \\ = & \{ (10) \text{ with common prefix } a \cdot \dagger \} \\ & u \sim *r \vee \dagger \cdot u \sim * \cdot a \cdot \dagger \cdot **r \end{aligned}$$

Now we are stuck because of the symbol a in the second disjunct.

If we could somehow eliminate this symbol, however, the road would be free to exploit our induction hypothesis $\dagger \cdot u \sim * \cdot \dagger \cdot **r \Rightarrow u \sim *r$. What we need is something like the eclipse lemma (11), except that the symbol that disappears is now a non-star:

$$(12) \quad u \sim * \cdot b \cdot r \Rightarrow u \sim *r \quad (\text{black-hole lemma})$$

Pictorially speaking, the star "swallows" an adjacent symbol, but this time there is no hope that the \Leftarrow direction holds as well. Again we start with $u := \varepsilon$:

$$\begin{aligned} & \varepsilon \sim * \cdot b \cdot r \\ = & \{ (7) \text{ with } p := b \cdot r \} \\ & \varepsilon \sim b \cdot r \\ = & \{ (5) \text{ with } p := r \} \\ & \text{false} \\ \Rightarrow & \{ \text{ex falso quodlibet} \} \\ & \varepsilon \sim *r \end{aligned}$$

Now the induction step $u := a \cdot u$:

$$\begin{aligned}
& a \cdot u \sim * \cdot b \cdot r \\
= & \{ (8) \text{ with } s, p := u, b \cdot r \} \\
& a \cdot u \sim b \cdot r \vee u \sim * \cdot b \cdot r \\
\Rightarrow & \{ \text{induction hypothesis} \} \\
& a \cdot u \sim b \cdot r \vee u \sim * \cdot r \\
= & \{ (6) \text{ with } s, p := u, r \} \\
& (a = b \wedge u \sim r) \vee u \sim * \cdot r \\
\Rightarrow & \{ \text{drop conjunct} \} \\
& u \sim r \vee u \sim * \cdot r \\
= & \{ (1) \text{ with } s := a \} \\
& (a \sim * \wedge u \sim r) \vee (a \sim * \wedge u \sim * \cdot r) \\
\Rightarrow & \{ (9) \text{ twice} \} \\
& a \cdot u \sim * \cdot r \vee a \cdot u \sim * \cdot * \cdot r \\
= & \{ (11) \} \\
& a \cdot u \sim * \cdot r \vee a \cdot u \sim * \cdot r \\
= & \{ \vee \text{ idempotent} \} \\
& a \cdot u \sim * \cdot r
\end{aligned}$$

This completes the proof of the black-hole lemma (12),
and we can resume the induction step $a \cdot t \cdot u \sim * \cdot a \cdot t \cdot * \cdot r \Rightarrow u \sim * \cdot r$.

After applying (8) and (10) to $a \cdot t \cdot u \sim * \cdot a \cdot t \cdot * \cdot r$, we continue:

$$\begin{aligned}
& u \sim * \cdot r \vee t \cdot u \sim * \cdot a \cdot t \cdot * \cdot r \\
\Rightarrow & \{ (12) \text{ with } u, b, r := t \cdot u, a, t \cdot * \cdot r \} \\
& u \sim * \cdot r \vee t \cdot u \sim * \cdot t \cdot * \cdot r \\
\Rightarrow & \{ \text{induction hypothesis} \} \\
& u \sim * \cdot r \vee u \sim * \cdot r \\
= & \{ \vee \text{ idempotent} \} \\
& u \sim * \cdot r
\end{aligned}$$

In summary, we have shown the easy part $u \sim * \cdot r \vdash t \cdot u \sim * \cdot t \cdot * \cdot r$ as well as the more challenging part $t \cdot u \sim * \cdot t \cdot * \cdot r \Rightarrow u \sim * \cdot r$. Thus, we can use the following theorem in the design of FINISH:

$$(13) \quad t \cdot u \sim * \cdot t \cdot * \cdot r \equiv u \sim * \cdot r \quad (\text{eclipse theorem})$$

Let us now analyze the remaining cases permitted by assertion $J \wedge (u = \varepsilon \vee q = \varepsilon \vee H \cdot u \neq H \cdot q)$.
If u and q are both empty, we have $t \sim * \cdot t$, which is always true:

$$\begin{aligned}
& \vdash \{ (0) \} \\
& \quad t \sim t \\
& \vdash \{ (1) \} \\
& \quad \varepsilon \sim *, \quad t \sim t \\
& \vdash \{ (2) \} \\
& \quad t \sim * \cdot t
\end{aligned}$$

If both are non-empty, let $u = a \cdot v$ and $q = b \cdot r$ with $a \neq b$:

$$\begin{aligned}
& t \cdot a \cdot v \sim * \cdot t \cdot b \cdot r \\
= & \{ (8) \} \\
& t \cdot a \cdot v \sim t \cdot b \cdot r \vee T.(t \cdot a \cdot v) \sim * \cdot t \cdot b \cdot r \\
= & \{ (10) \} \\
& a \cdot v \sim b \cdot r \vee T.(t \cdot a \cdot v) \sim * \cdot t \cdot b \cdot r \\
= & \{ (6) \text{ with } a \neq b \} \\
& \text{false} \vee T.(t \cdot a \cdot v) \sim * \cdot t \cdot b \cdot r \\
= & \{ \text{drop false disjunct} \} \\
& T.(t \cdot a \cdot v) \sim * \cdot t \cdot b \cdot r
\end{aligned}$$

If only u is empty, we have $t \sim * \cdot t \cdot b \cdot r$ and see intuitively that the pattern is "too long", yielding the proof obligation:

$$(14) \ t \sim * \cdot t \cdot b \cdot r \equiv \text{false}$$

And if only q is empty, $t \cdot a \cdot v \sim * \cdot t$ means that $t \cdot a \cdot v$ ends with suffix t , so we can immediately shrink the sample to the appropriate length.

Even without a formal treatment of the last two cases, the design of FINISH is already clear from $t \sim * \cdot t$, (14), $t \cdot a \cdot v \sim * \cdot t \cdot b \cdot r \equiv T.(t \cdot a \cdot v) \sim * \cdot t \cdot b \cdot r$, and the eclipse theorem (13).

FINISH:

```

{ J  $\wedge$  (u =  $\epsilon$   $\vee$  q =  $\epsilon$   $\vee$  H.u  $\neq$  H.q) }
if q =  $\epsilon$   $\rightarrow$ 
  if u =  $\epsilon$   $\rightarrow$  y := true
  [] u  $\neq$   $\epsilon$   $\rightarrow$  SHRINK
  fi
[] q  $\neq$   $\epsilon$   $\rightarrow$ 
  if H.q  $\neq$  *  $\rightarrow$ 
    if u =  $\epsilon$   $\rightarrow$  n := true
    [] u  $\neq$   $\epsilon$   $\rightarrow$  s := T.s
    fi
  [] H.q = *  $\rightarrow$  s, p := u, q
  fi
fi
{ I }

```

With regard to termination, we notice that each branch achieves progress by setting a flag or shortening a string. In the last branch, the sample might stay the same (if $t = \epsilon$), but the pattern always loses its first star.

Now to the outstanding proof of (14), which calls for a counting argument over the non-star symbols in a pattern. Guided by calculus (0 - 2), we define a "non-star counting function" N from patterns to natural numbers:

$$(15) \ N.s = |s|$$

$$(16) \ N.* = 0$$

$$(17) \ N.(p \cdot q) = N.p + N.q$$

The expression $|s|$ above denotes the length of s measured in symbols.
 As a next step, we postulate the following relationship between \sim and N :

$$(18) s \sim p \Rightarrow |s| \geq N.p \quad (\text{general counting lemma})$$

The proof is by induction for all wildcard matches enumerated by calculus (0 - 2).

Case 0: $\vdash s \sim s$

$$\begin{aligned} & |s| \geq N.s \\ = & \{ (15) \} \\ & |s| \geq |s| \\ = & \{ \geq \text{reflexive} \} \\ & \text{true} \end{aligned}$$

Case 1: $\vdash s \sim *$

$$\begin{aligned} & |s| \geq N.* \\ = & \{ (16) \} \\ & |s| \geq 0 \\ = & \{ \text{string theory} \} \\ & \text{true} \end{aligned}$$

Case 2: $s \sim p, t \sim q \vdash s \cdot t \sim p \cdot q$

$$\begin{aligned} & |s| \geq N.p \wedge |t| \geq N.q \\ \Rightarrow & \{ \text{arithmetic} \} \\ & |s| + |t| \geq N.p + N.q \\ = & \{ \text{string theory} \} \\ & |s \cdot t| \geq N.p + N.q \\ = & \{ (17) \} \\ & |s \cdot t| \geq N.(p \cdot q) \end{aligned}$$

So if the pattern $* \cdot t \cdot b \cdot r$ in (14) is really "too long" for sample t ,
 we should now be able to write this down as a formal calculation:

$$\begin{aligned} & N.(* \cdot t \cdot b \cdot r) \\ = & \{ (17) \text{ three times} \} \\ & N.* + N.t + N.b + N.r \\ = & \{ (16) ; (15) \text{ twice} \} \\ & 0 + |t| + 1 + N.r \\ > & \{ \text{arithmetic} \} \\ & |t| \end{aligned}$$

The result looks like a negated RHS of (18) (with $s, p := t, * \cdot t \cdot b \cdot r$),
 from which the negation of the LHS follows: $|t| < N.(* \cdot t \cdot b \cdot r) \Rightarrow \neg t \sim * \cdot t \cdot b \cdot r$

This completes the proof of (14), and the only thing left to do is the design of SHRINK satisfying:

$$\begin{aligned} & \{ I \wedge s = t \cdot u \wedge u \neq \varepsilon \wedge p = * \cdot t \} \\ & \text{SHRINK} \\ & \{ I \} \end{aligned}$$

What follows is essentially another counting argument, but this time without any stars in the pattern. The impact is that the number of non-stars $N.p$ is simply the pattern length $|p|$, and the "z" from the general counting lemma (18) changes to equality:

$$(19) s \sim t \Rightarrow |s| = |t| \quad (\text{special counting lemma})$$

The proof is analogous to the proof of (18), except that we omit axiom (1).

Case 0: $\vdash s \sim s$

$$\begin{aligned} & |s| = |s| \\ = & \{ = \text{ reflexive } \} \\ & \text{true} \end{aligned}$$

Case 2: $s \sim p, t \sim q \vdash s \cdot t \sim p \cdot q$

$$\begin{aligned} & |s| = |p| \wedge |t| = |q| \\ \Rightarrow & \{ \text{Leibniz' Principle} \} \quad \underline{\text{Remark}} \text{ I.e. function application preserves equality. (End of Remark.)} \\ & |s| + |t| = |p| + |q| \\ = & \{ \text{string theory} \} \\ & |s \cdot t| = |p \cdot q| \end{aligned}$$

Applied to the design of SHRINK, this means that as long as $|s| > |t|$, the first disjunct of

$$s \sim * \cdot t \equiv s \sim t \vee \top.s \sim * \cdot t$$

is always false, so $s := \top.s$ preserves \mathbf{I} . The precondition of SHRINK even tells us how often we can do this, because $s = t \cdot u$ means that the length difference to be eliminated equals $|u|$.

Well, and when $|s| = |t|$, either s is empty or $\top.s \sim * \cdot t$ violates the general counting lemma (18) because $|\top.s| = |s| - 1$ and $N.(* \cdot t) = 0 + |t|$. In both cases, $s \sim * \cdot t \equiv s \sim t$.

The simplest way to iterate $|u|$ times is to shrink s and u simultaneously and make $u \neq \varepsilon$ the guard of the loop. Another benefit is the invariance of $|s| = |t \cdot u|$, the "counting version" of $s = t \cdot u$:

$$\begin{aligned} & \text{wp} . (s, u := \top.s, \top.u) . ((M \equiv s \sim * \cdot t) \wedge |s| = |t \cdot u|) \\ = & \{ \text{axiom of assignment} \} \\ & (M \equiv \top.s \sim * \cdot t) \wedge |\top.s| = |t \cdot \top.u| \\ = & \{ \text{string theory} \} \\ & (M \equiv \top.s \sim * \cdot t) \wedge |s| - 1 = |t \cdot u| - 1 \wedge u \neq \varepsilon \\ = & \{ \text{arithmetic} \} \\ & (M \equiv \top.s \sim * \cdot t) \wedge |s| = |t \cdot u| \wedge u \neq \varepsilon \\ = & \{ (19) \text{ considering } |s| > |t| \} \\ & (M \equiv s \sim t \vee \top.s \sim * \cdot t) \wedge |s| = |t \cdot u| \wedge u \neq \varepsilon \\ = & \{ (8) \text{ with } a, s, p := \top.s, \top.s, t \} \\ & (M \equiv s \sim * \cdot t) \wedge |s| = |t \cdot u| \wedge u \neq \varepsilon \end{aligned}$$

Now we have all the ingredients to spell out the missing part of FINISH.

```

SHRINK:
  { I ∧ s = t·u ∧ u ≠ ε ∧ p = *·t }
L: { I ∧ p = *·t ∧ |s| = |t·u| }
  do u ≠ ε →
    { L ∧ u ≠ ε }
    s, u := T.s, T.u
    { L }
  od;
  { L ∧ u = ε }
  { I ∧ p = *·t ∧ |s| = |t| }
  p := T.p
  { I }

```

So not only does SHRINK re-establish **I** (for which a simple $s := T.s$ would have sufficed), but it also optimizes away subsequent PUSH;FINISH cycles doomed to failure because of left-over symbols in one string. This completes the design of our program.

* * *

A few final remarks are in order:

(i) The program variable **t** (capturing the common prefix in FRAG8) is vital for the correctness proof but largely ignored by the program itself - except for the two maintenance assignments $t := ε$ and $t := t·H.u$. Hence, for execution on a real machine, the variable and both assignments can be optimized away.

(ii) Another good candidate for an optimization is the point between $q := T.p$ and PUSH in FRAG8. If q is empty there, PUSH has nothing to do, but it may take a while to SHRINK away the remaining sample. On the other hand, axiom (1) leaves no doubt that the final answer is "yes", so we may as well immediately set the flag **y**.

```

FRAG8:
  { I ∧ s ≠ ε ∧ H.p = * }
  u, q := s, T.p;
  { I ∧ p = *·q }
  if q = ε →
    { I ∧ p = * }
    { I ∧ M }
    y := true
  [] q ≠ ε →
    PUSH;
    FINISH
  fi
  { I }

```

(iii) Our initial assumption of star-free samples (formally captured by $* \notin \Sigma$) may be nice in theory but hard to enforce in practice. So let us postulate the existence of two "different" star symbols, together with the convention that a sample star never equals a pattern star. In fact, most of our program implicitly works that way due to special guards for pattern stars. However, the plain comparison $H.u = H.q$ in the PUSH loop needs the extra conjunct $H.q \neq *$ to enforce inequality among two stars:

```
if H.u = H.q  $\wedge$  H.q  $\neq$  *  $\rightarrow$  u, q := T.u, T.q
[] H.u  $\neq$  H.q  $\vee$  H.q = *  $\rightarrow$  d := true
fi
```

(iv) After elimination of \dagger (which grows from ϵ to a new string of appended symbols), only the following string instructions remain in the program:

- a) Check if string = ϵ
- b) Extract H.string
- c) Assign T.string to a variable

Thus, even the most primitive string representations - like pointers to zero-terminated memory chunks - are graciously accommodated:

- a) Compare pointer target with zero
- b) Read symbol from pointer target
- c) Store incremented pointer value

(v) Despite our careful design of FINISH and SHRINK, there is still room to reduce the worst-case asymptotic running time if we resort to highly sophisticated string-matching algorithms, like the linear-time Knuth-Morris-Pratt solution. However, since our goal was an orderly design from first principles (without pulling any rabbits), such advanced optimizations exceed the scope of this work.

Acknowledgments

Although I never met Dijkstra in person, I got a faithful exposure to his ideas thanks to Edgar Knapp, who had come from Austin to Purdue in 1992, where I became his student one year later. Credit goes to my colleague Gunter Hick, whose lunchtime remark on a non-recursive four-pointer wildcard algorithm triggered my interest in the topic - and again to Edgar Knapp, as an inspiring brainstorming and review partner for its calculational treatment. (End of Acknowledgments.)

APPENDIX: The complete program expressed in "Guarded C"

```

#include <stdio.h>
#include <stdbool.h>
#include <assert.h>

#define SKIP {}
#define ABORT assert (false)

#define IF if (
#define FI } else ABORT;

#define DO for (;;) { if (
#define OD } else break; }

#define THEN ) {
#define II } else if (

bool empty (char* string) { return *string == 0; }
char head (char* string) { return *string; }
char* tail (char* string) { return string + 1; }

main (int argc, char* argv [])
{
    char* s = argv [1];
    char* p = argv [2];

    bool y = false;
    bool n = false;

    DO y == n THEN
        IF empty (p) THEN
            IF empty (s) THEN y = true;
            II!empty (s) THEN n = true;
            FI
        II!empty (p) THEN
            IF head (p) != '*' THEN
                IF empty (s) THEN n = true;
                II!empty (s) THEN
                    IF head (s) == head (p) THEN s = tail (s); p = tail (p);
                    II head (s) != head (p) THEN n = true;
                    FI
                FI
            II head (p) == '*' THEN
                IF empty (s) THEN p = tail (p);
                II!empty (s) THEN
                    char* u = s; char* q = tail (p);
                    IF empty (q) THEN y = true;
                    II!empty (q) THEN
                        bool d = false;
                        DO !empty (u) && !empty (q) && !d THEN
                            IF head (u) == head (q) && head (q) != '*' THEN
                                u = tail (u); q = tail (q);
                                II head (u) != head (q) || head (q) == '*' THEN
                                    d = true;
                                FI
                            OD
                        IF empty (q) THEN
                            IF empty (u) THEN y = true;
                            II!empty (u) THEN
                                DO !empty (u) THEN
                                    s = tail (s); u = tail (u);
                                OD
                                p = tail (p);
                            FI
                        II!empty (q) THEN
                            IF head (q) != '*' THEN
                                IF empty (u) THEN n = true;
                                II!empty (u) THEN s = tail (s);
                                FI
                            II head (q) == '*' THEN s = u; p = q;
                            FI
                        FI
                    FI
                FI
            FI
        FI
    OD
    if (y) printf ("yes\n");
    if (n) printf ("no\n");
}

```