

We use CFGs to describe both human languages and programming languages.

So we want to be able to compute whether a given string is in the language of a particular CFG:

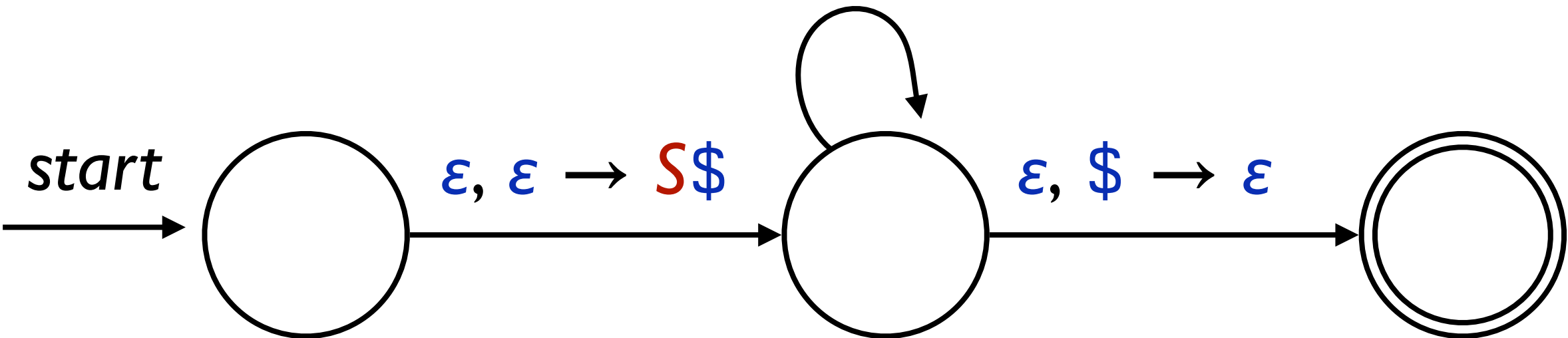
Is this a well-formed utterance in English?

Is this a syntactically correct Python program?

Why isn't this good enough?

$S \rightarrow 1S1$
 $S \rightarrow 1S$
 $S \rightarrow \geq$

$\epsilon, S \rightarrow 1S1$
 $\epsilon, S \rightarrow 1S$
 $\epsilon, S \rightarrow \geq$
 $\Sigma, \Sigma \rightarrow \epsilon$



Given a context-free grammar G and a word w , we already have a way to tell if G generates w ...

...using nondeterministic guessing of the right rule to use at each step.

We need to be able to do this *algorithmically* – and we'd like it to be fast!

Two parts:

1. Convert the grammar to a “normal form”
2. Use dynamic programming to efficiently parse the input

The definition of a context-free grammar permits unlimited flexibility in the form of the rule bodies.

Advantage: Flexibility in *designing* grammars

Disadvantage: Difficulty in *analyzing* grammars

A *normal form* grammar obeys certain restrictions on the form of the rules.

It can simplify derivations.

It makes proofs about grammars easier because there are fewer possibilities to consider.

It eliminates redundancy.

Properties of normal forms:

Every CFL can be generated by a normal form grammar.

There are algorithmic transformations to convert any CFG into an equivalent grammar in normal form.

Chomsky normal form (CNF)

Each rule is one of these forms:

$A \rightarrow BC$, where neither B nor C is S

$A \rightarrow a$

$S \rightarrow \epsilon$

Chomsky normal form (CNF)

Each rule is one of these forms:

$A \rightarrow BC$, where neither B nor C is S

$A \rightarrow a$

$S \rightarrow \varepsilon$

Note: Only the start symbol can produce ε

THEOREM Any context-free language is generated by a context-free grammar in Chomsky normal form.

PROOF IDEA Show how to convert a CFG into CNF.

Conversion requires a sequence of rule additions, deletions, and modifications that leave the set of strings generated by the grammar unchanged.

*The details are in the reading – and in an appendix to these slides.
We'll just walk through an example in class.*

Example conversion

$$T \rightarrow aTb \mid \varepsilon$$

1. *Make start symbol non-recursive*
2. *Eliminate all ε -rules (except $S \rightarrow \varepsilon$)*
3. *Eliminate all chain/unit rules*
4. *Eliminate useless symbols*
5. *Split rules if necessary*

Example conversion

$$S \rightarrow T$$

$$T \rightarrow aTb \mid \varepsilon$$

1. *Make start symbol non-recursive*
2. *Eliminate all ε -rules (except $S \rightarrow \varepsilon$)*
3. *Eliminate all chain/unit rules*
4. *Eliminate useless symbols*
5. *Split rules if necessary*

Example conversion

$$S \rightarrow T \mid \varepsilon$$

$$T \rightarrow aTb \mid ab$$

1. *Make start symbol non-recursive*
2. **Eliminate all ε -rules (except $S \rightarrow \varepsilon$)**
3. *Eliminate all chain/unit rules*
4. *Eliminate useless symbols*
5. *Split rules if necessary*

Example conversion

$S \rightarrow aTb \mid ab \mid \varepsilon$

$T \rightarrow aTb \mid ab$

1. *Make start symbol non-recursive*
2. *Eliminate all ε -rules (except $S \rightarrow \varepsilon$)*
3. ***Eliminate all chain/unit rules***
4. *Eliminate useless symbols*
5. *Split rules if necessary*

Example conversion

$S \rightarrow aTb \mid ab \mid \varepsilon$

$T \rightarrow aTb \mid ab$

1. *Make start symbol non-recursive*
2. *Eliminate all ε -rules (except $S \rightarrow \varepsilon$)*
3. *Eliminate all chain/unit rules*
4. ***Eliminate useless symbols***
5. *Split rules if necessary*

None!

Example conversion

$$S \rightarrow V_a T V_b \mid V_a V_b \mid \varepsilon$$

$$T \rightarrow V_a T V_b \mid V_a V_b$$

$$V_a \rightarrow a$$

$$V_b \rightarrow b$$

1. Make start symbol non-recursive
2. Eliminate all ε -rules (except $S \rightarrow \varepsilon$)
3. Eliminate all chain/unit rules
4. Eliminate useless symbols
- 5. Split rules if necessary**

a. Add new rules for terminals

Example conversion

$$S \rightarrow V_a C \mid V_a V_b \mid \epsilon$$

$$T \rightarrow V_a C \mid V_a V_b$$

$$C \rightarrow T V_b$$

$$V_a \rightarrow a$$

$$V_b \rightarrow b$$

1. Make start symbol non-recursive
2. Eliminate all ϵ -rules (except $S \rightarrow \epsilon$)
3. Eliminate all chain/unit rules
4. Eliminate useless symbols
- 5. Split rules if necessary**

b. Replace rules with three or more symbols on the right side by creating new rules with new variables

And now the grammar is in
Chomsky normal form.



*I approve this
grammar!*

CYK parsing algorithm

PROBLEM Given a CFG $G = (V, \Sigma, P, S)$ and a string $w \in \Sigma^*$, determine if w is in $L(G)$, that is, whether $S \xRightarrow{*} w$.

The Cocke–Younger–Kasami (CYK) algorithm is an $O(n^3)$ algorithm ($n = \text{length of input string}$) that uses dynamic programming to determine the derivability of the string from the grammar.

PROBLEM Given a CFG $G = (V, \Sigma, P, S)$ and a string $w \in \Sigma^*$, determine if w is in $L(G)$, that is, whether $S \xRightarrow{*} w$.

The Cocke–Younger–Kasami (CYK) algorithm is an $O(n^3)$ algorithm ($n = \text{length of input string}$) that uses dynamic programming to determine the derivability of the string from the grammar.

If you haven't seen "Big O Notation" before, it means that it takes at most n^3 steps of computation (loosely defined) to process an input of length n .

PROBLEM Given a CFG $G = (V, \Sigma, P, S)$ and a string $w \in \Sigma^*$, determine if w is in $L(G)$, that is, whether $S \xRightarrow{*} w$.

The Cocke–Younger–Kasami (CYK) algorithm is an $O(n^3)$ algorithm ($n = \text{length of input string}$) that **uses dynamic programming** to determine the membership of the string from the grammar.

Dynamic programming is a class of methods that avoid duplicate computation at the expense of memory.

Values that may be used in future computations are stored in a table.

PROBLEM Given a CFG $G = (V, \Sigma, P, S)$ and a string $w \in \Sigma^*$, determine if w is in $L(G)$, that is, whether $S \xRightarrow{*} w$.

The Cocke–Younger–Kasami (CYK) algorithm is an $O(n^3)$ algorithm ($n = \text{length of input string}$) that uses dynamic programming to determine the derivability of the string from the grammar.

CYK algorithm

Start with a Chomsky normal form grammar for L .

Build a two-dimensional table to keep track of which portions of the input string can be generated starting from which variables.

Table

The horizontal axis corresponds to the positions of the string

$$w = x_1x_2\dots x_n.$$

Table entry $X_{i,j}$ is the set of non-terminals A such that

$$A \xRightarrow{*} w_i w_{i+1} \dots w_j.$$

We are particularly interested in whether

S is in $X_{1,n}$ because that is the same as

saying $S \xRightarrow{*} w$ (that is, w is in L)

| | | | | | |
|-----|-----|-----|-----|-----|---|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| | a | a | b | b | b |

CYK algorithm strategy

Given a string $w = w_1 w_2 w_3 \dots w_n$,
let $w_{i,j}$ be substring $w_i \dots w_j$ of w .

1. For each substring $w_{i,j}$ of *length 1*, find set $X_{i,j}$ of all variables with rule $A \rightarrow w_{i,j}$.
2. For each substring $w_{i,j+1}$ of *length 2*, find set $X_{i,j+1}$ of all variables that derive $w_{i,j+1}$

...

$n-1$. For substrings $w_{1,n-1}, w_{2,n}$ of *length $n-1$* , find sets $X_{1,n-1}, X_{2,n}$ of all variables that derive $w_{1,n-1}, w_{2,n}$.

n . For the string $w_{1,n} = w$ of *length n* , find the set $X_{1,n}$ of all variables that initiate derivation of w .

At the end, if $S \in X_{1,n}$, then
 $w \in L(G)$

Example

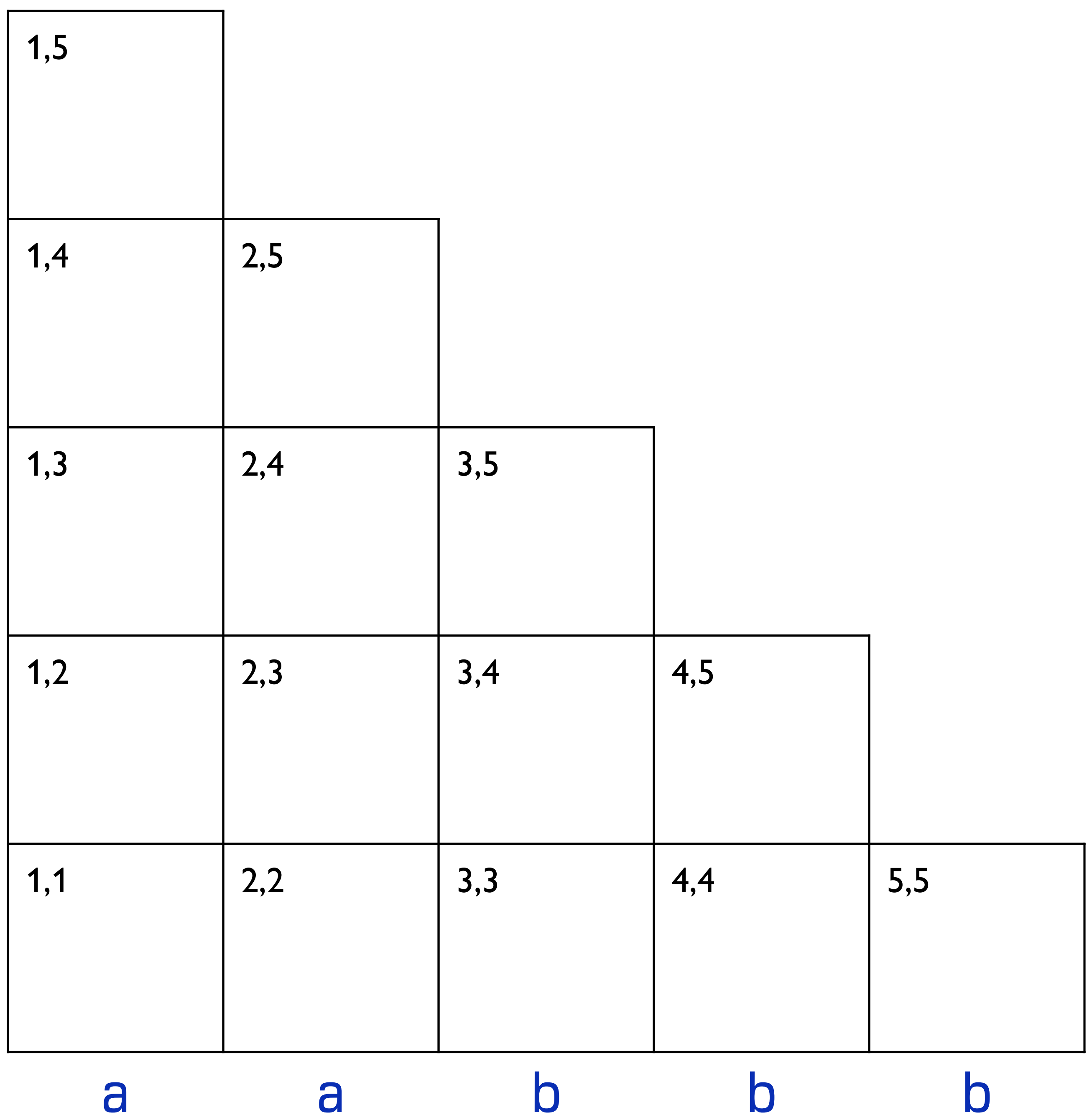
We'll use the algorithm to determine if the string $w = \text{aabbb}$ is in the language generated by the grammar

$$S \rightarrow AB$$

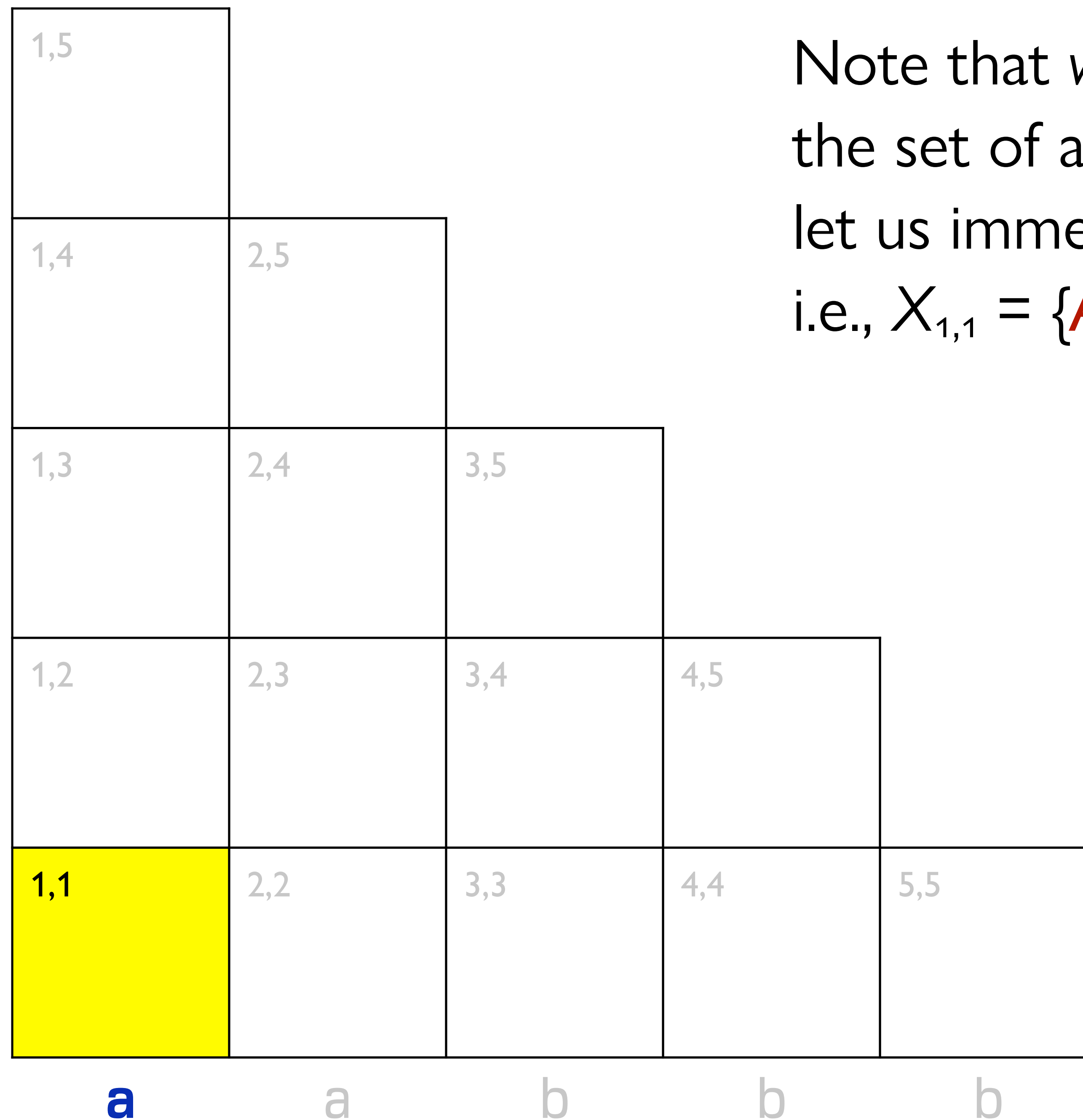
$$A \rightarrow BB \mid a$$

$$B \rightarrow AB \mid b$$

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

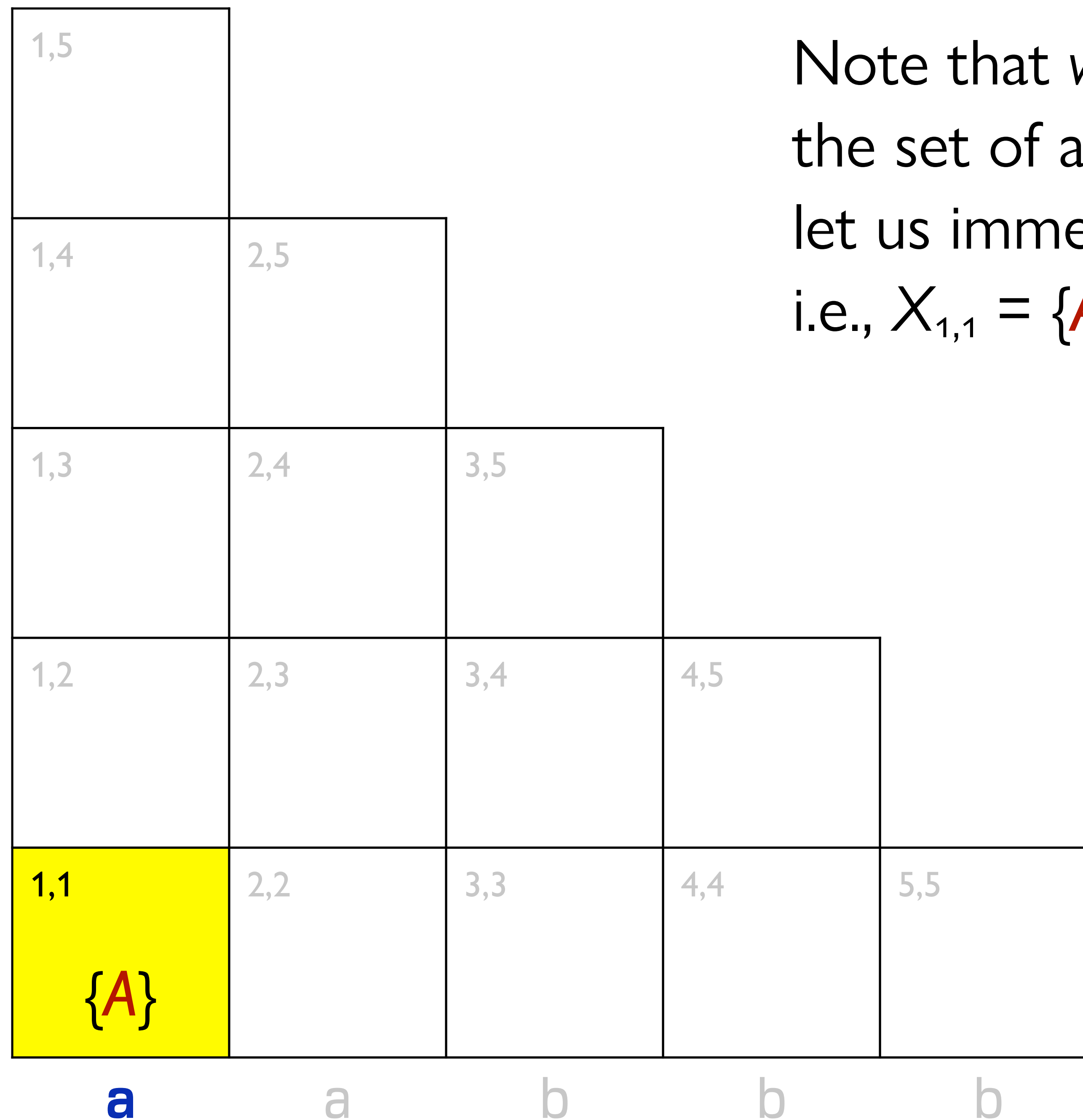


$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



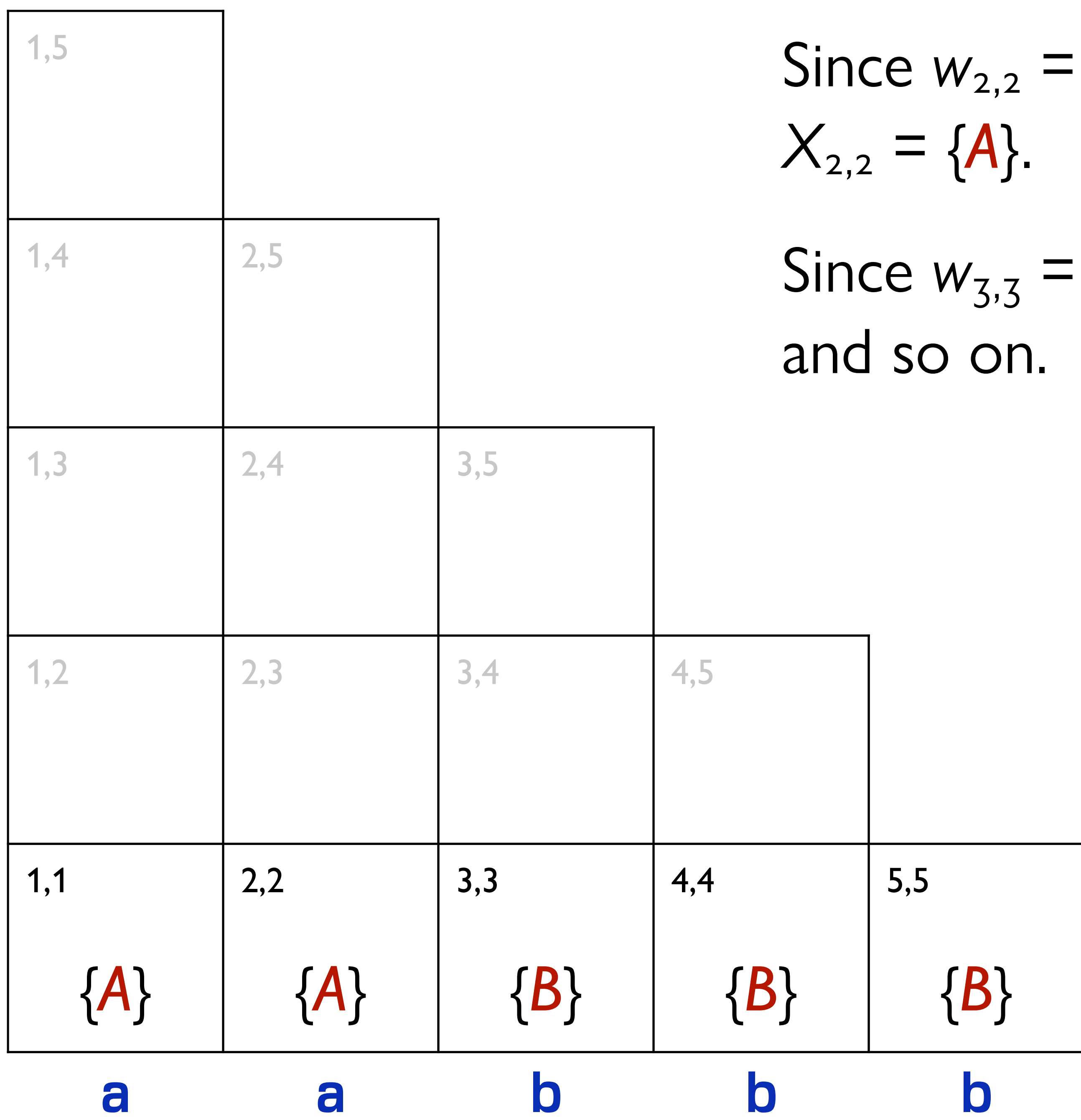
Note that $w_{1,1} = a$, so $X_{1,1}$ is the set of all variables that let us immediately derive a , i.e., $X_{1,1} = \{A\}$.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



Note that $w_{1,1} = a$, so $X_{1,1}$ is the set of all variables that let us immediately derive a , i.e., $X_{1,1} = \{A\}$.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



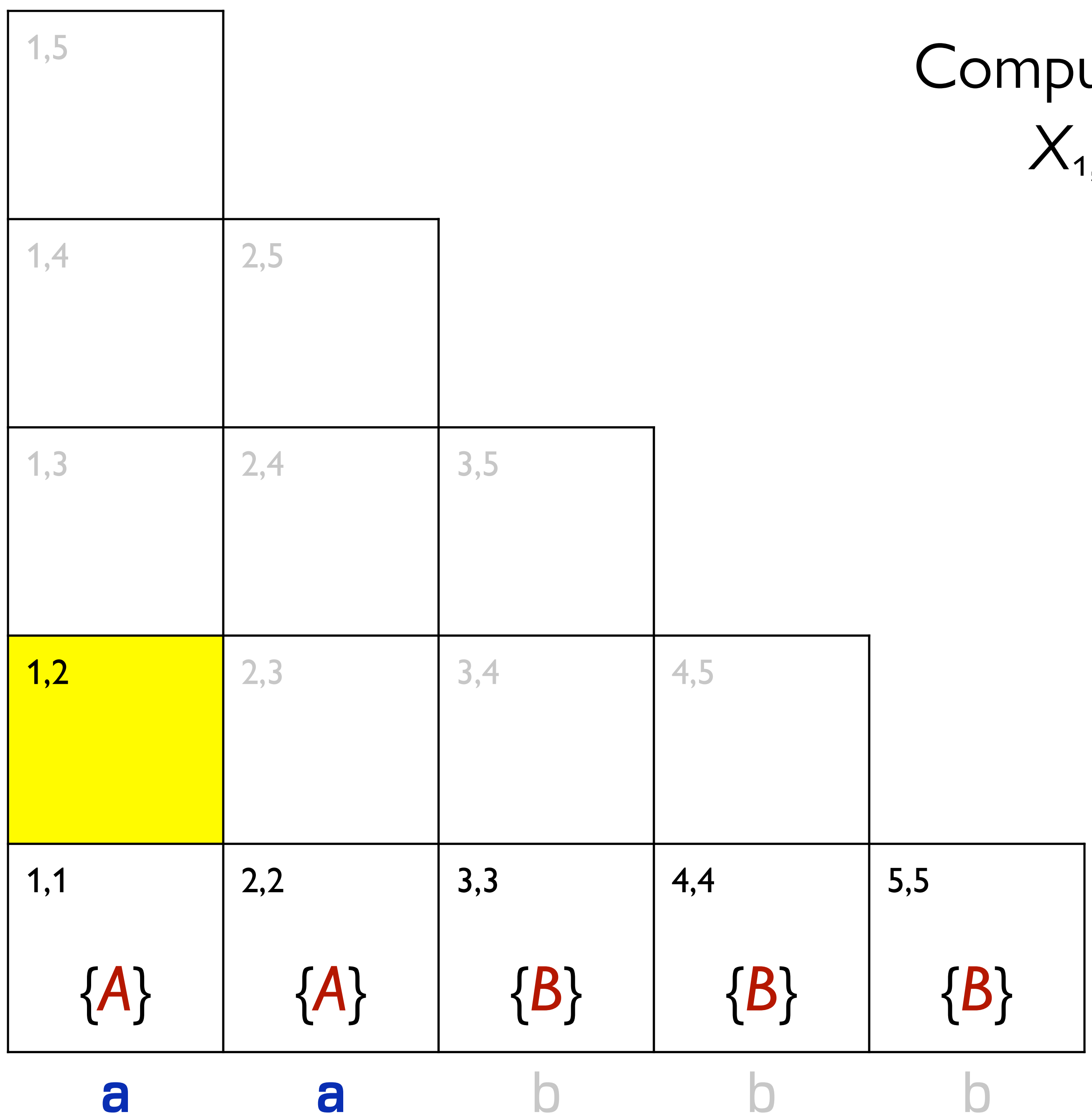
Since $w_{2,2} = a$, we also have $X_{2,2} = \{A\}$.

Since $w_{3,3} = b$, $X_{3,3} = \{B\}$, and so on.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

Compute

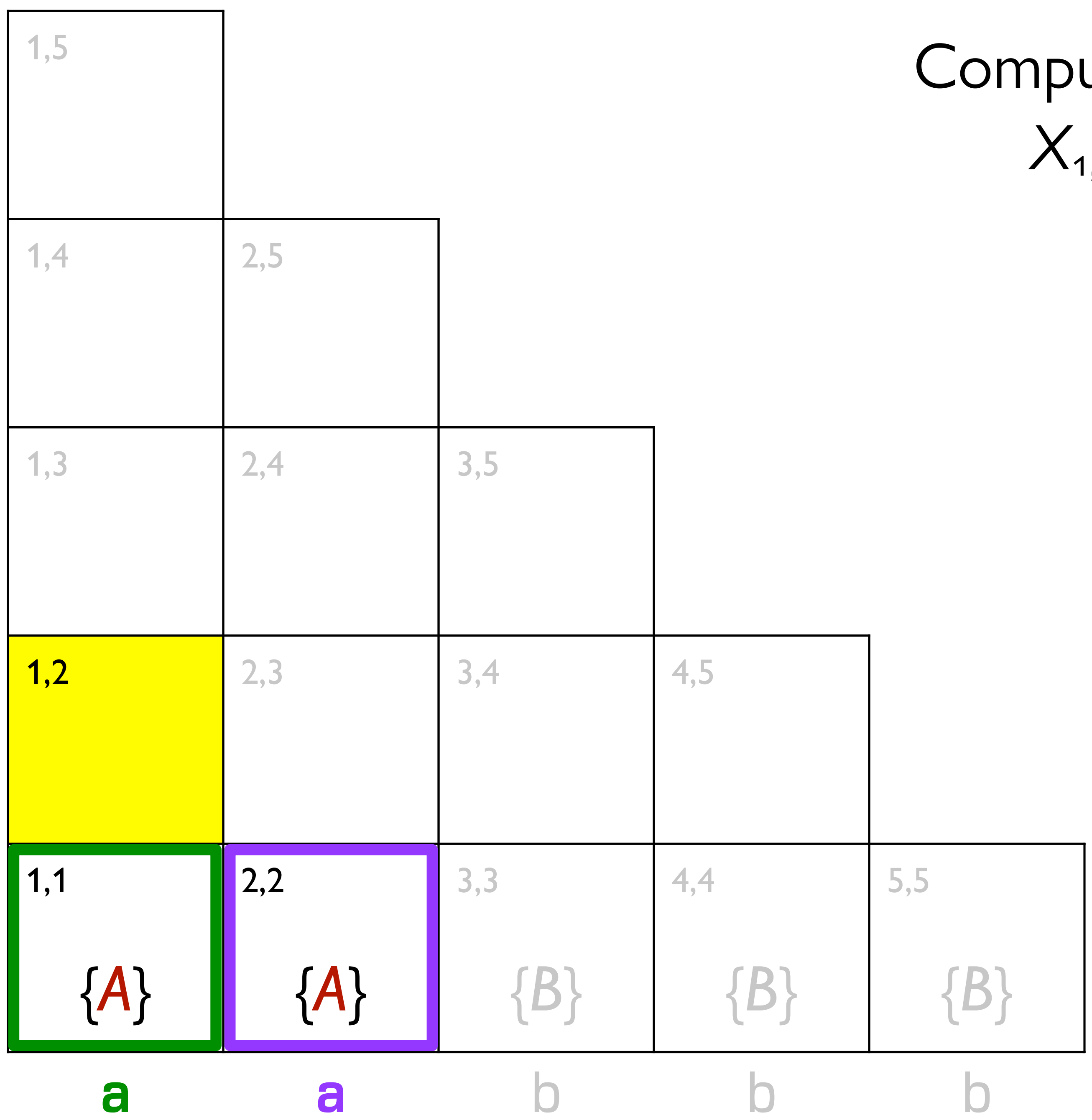
$$X_{1,2} = \{P \mid P \rightarrow QR, \\ Q \in X_{1,1}, \\ R \in X_{2,2}\}.$$



$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

Compute

$$X_{1,2} = \{P \mid P \rightarrow QR, Q \in X_{1,1}, R \in X_{2,2}\}.$$



$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|----------|----------|-----|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| | | | | | |
| a | a | b | b | b | |

Cell (1,2) contains \emptyset and is highlighted in yellow.
 Cell (1,1) contains $\{A\}$ and is highlighted with a green border.
 Cell (2,2) contains $\{A\}$ and is highlighted with a purple border.
 Cells (3,3), (4,4), and (5,5) contain $\{B\}$.

Compute

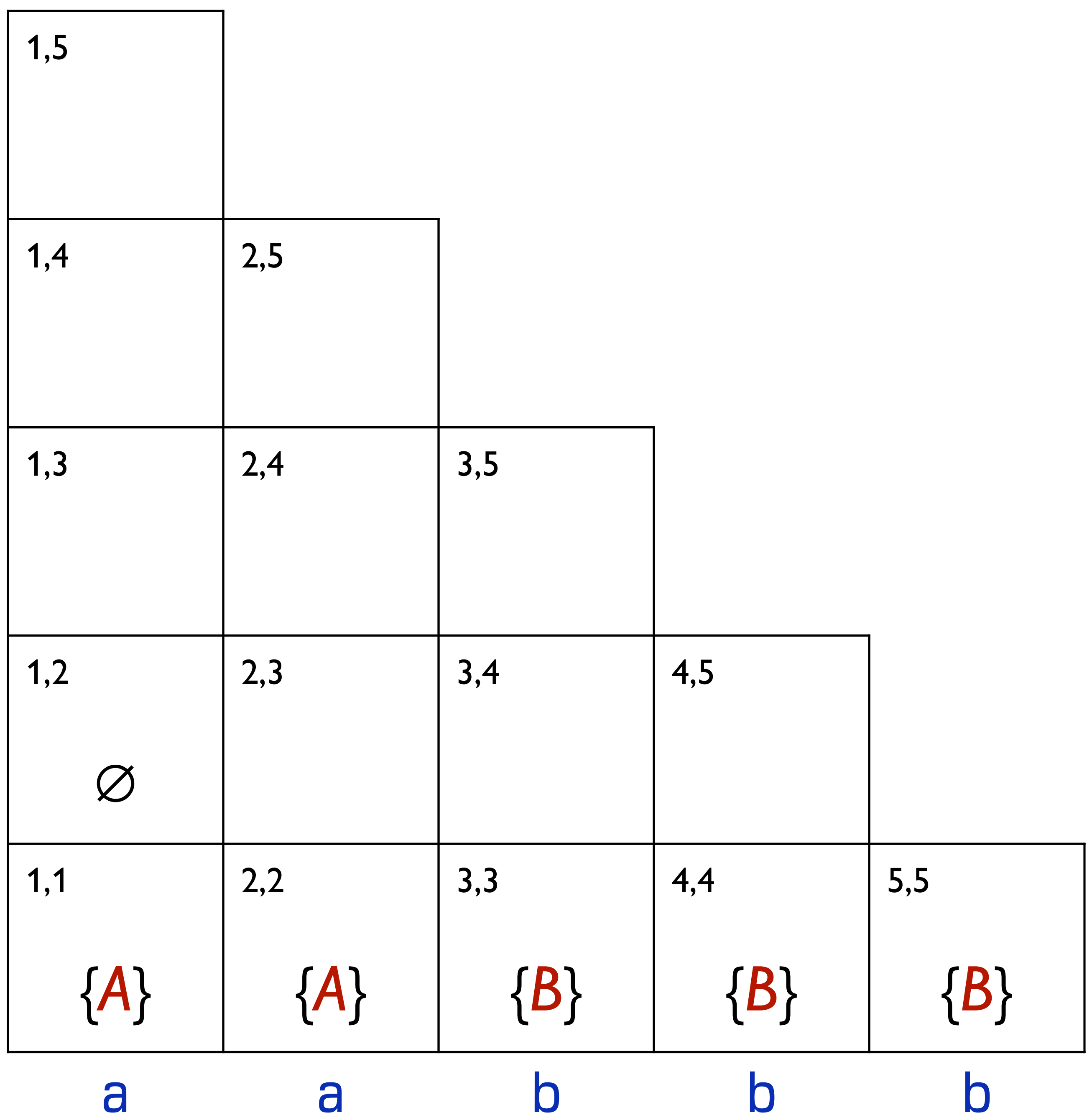
$$X_{1,2} = \{P \mid P \rightarrow QR, \\
 Q \in X_{1,1}, \\
 R \in X_{2,2}\}.$$

Since $X_{1,1} = \{A\}$ and

$$X_{2,2} = \{A\},$$

$X_{1,2} = \{P \mid P \rightarrow AA\}$, and
 there are no rules with
 this body.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



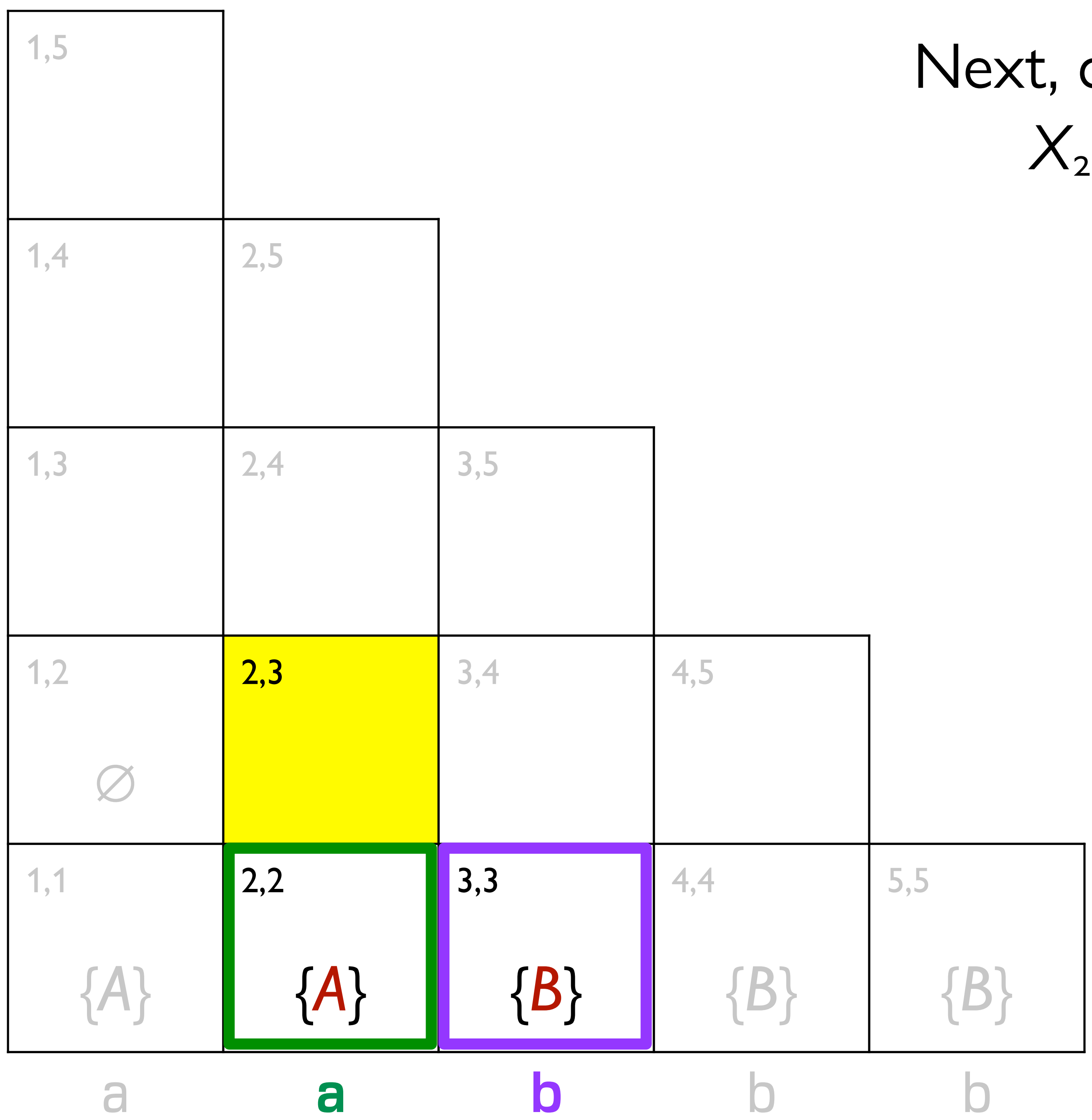
$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

Next, compute

$$X_{2,3} = \dots$$

| | | | | | |
|-----|-----|-----|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



Next, compute

$$X_{2,3} = \{P \mid P \rightarrow QR, Q \in X_{2,2}, R \in X_{3,3}\}.$$

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|-----|-----|-----|-----|---|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| | {A} | {B} | {B} | {B} | |
| | a | a | b | b | b |

Next, compute

$$X_{2,3} = \{P \mid P \rightarrow QR, Q \in X_{2,2}, R \in X_{3,3}\}.$$

The required right side is AB , thus $X_{2,3} = \{S, B\}$.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|-----|--------|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| a | a | b | b | b | |
| | {A} | {B} | {B} | {B} | |
| | ∅ | {S, B} | | | |

Next, compute

$$X_{2,3} = \{P \mid P \rightarrow QR, Q \in X_{2,2}, R \in X_{3,3}\}.$$

The required right side is AB , thus $X_{2,3} = \{S, B\}$.

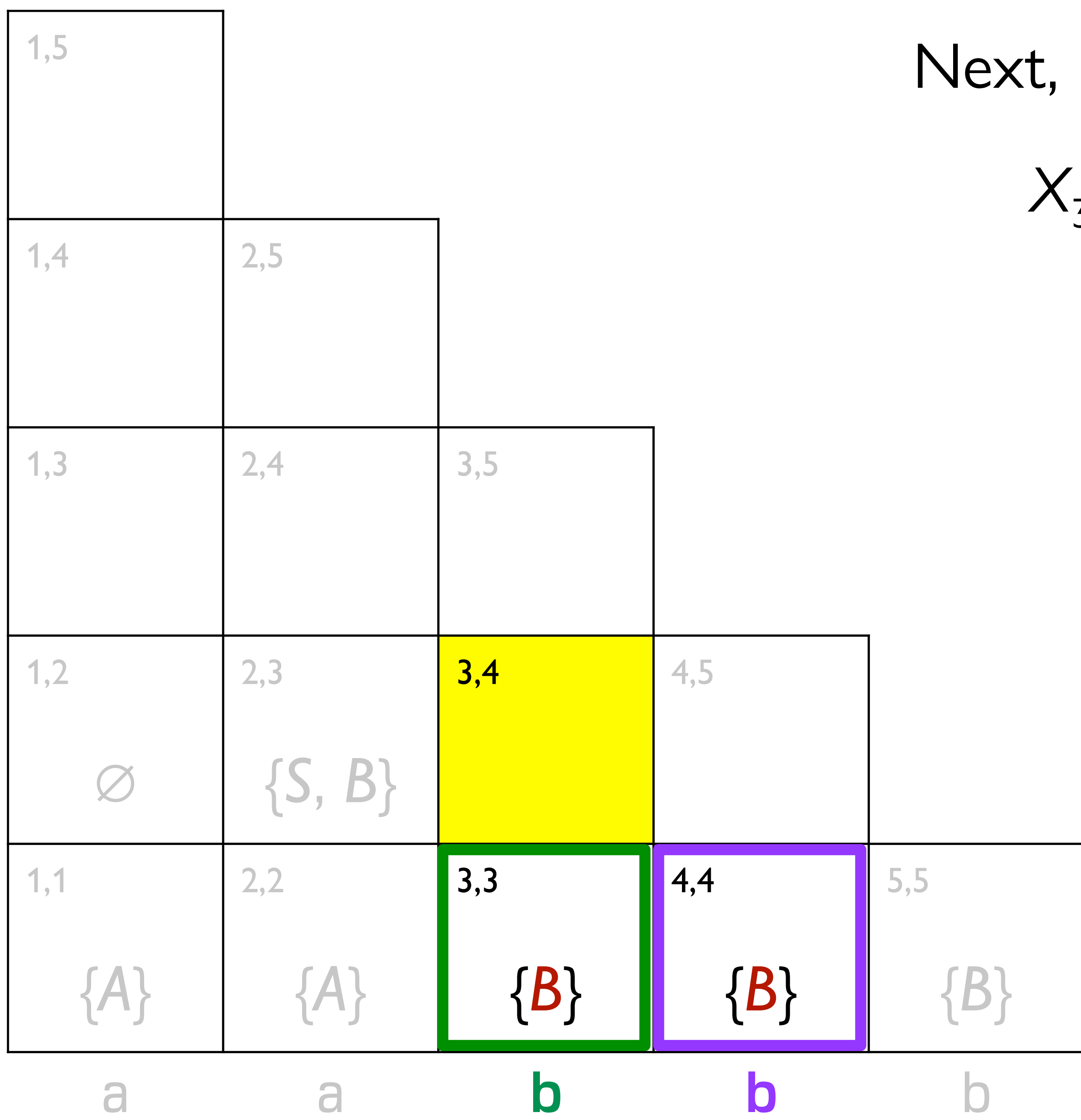
$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

Next,

$$X_{3,4} = \dots$$

| | | | | | |
|-----|-----|-----|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



Next,

$$X_{3,4} = \{P \mid P \rightarrow QR, Q \in X_{3,3}, R \in X_{4,4}\}$$

| |
|---------------------------|
| $S \rightarrow AB$ |
| $A \rightarrow BB \mid a$ |
| $B \rightarrow AB \mid b$ |

| | | | | |
|-----|-----|-----|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| a | a | b | b | b |

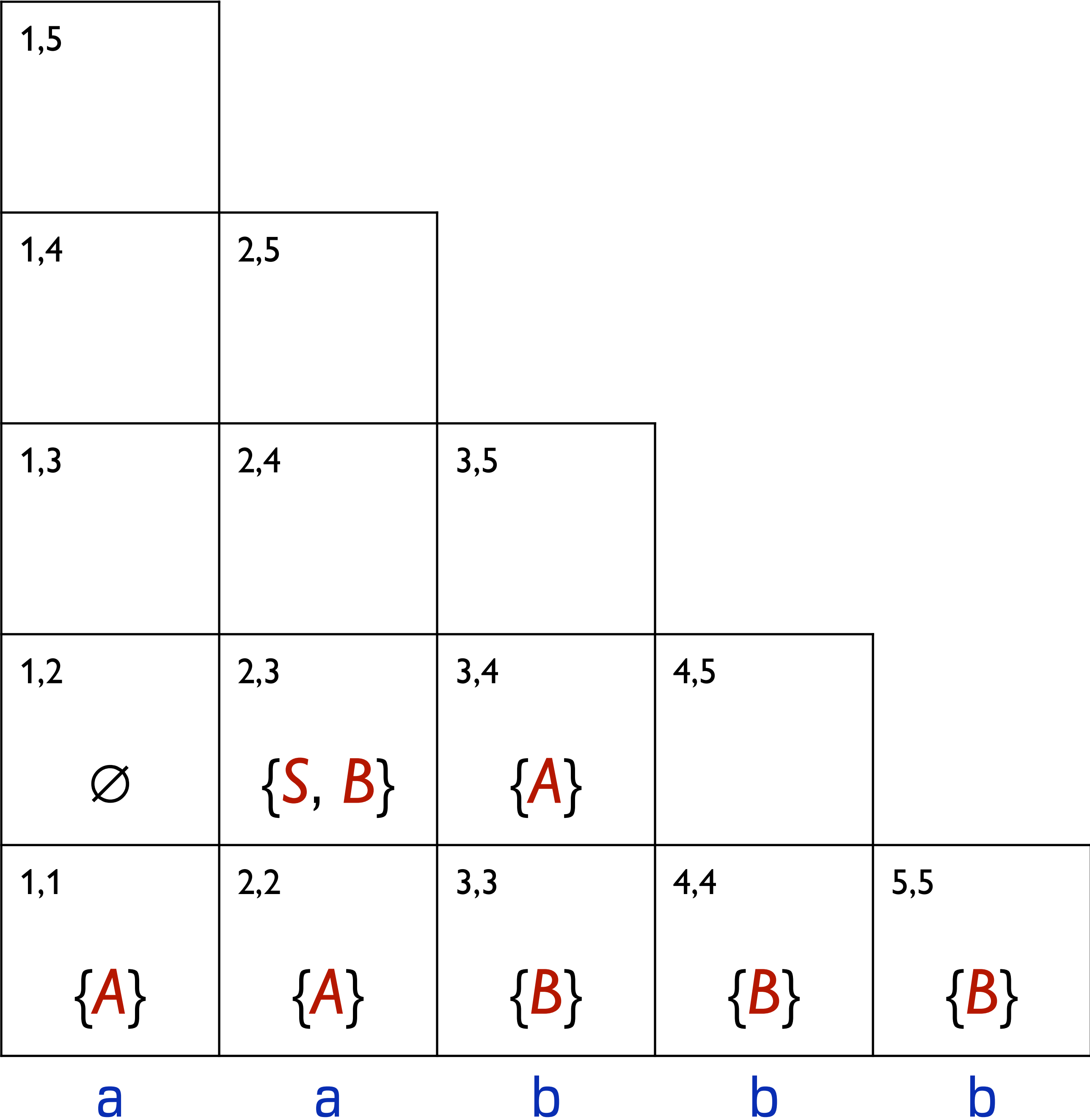
Next,

$$X_{3,4} = \{P \mid P \rightarrow QR, Q \in X_{3,3}, R \in X_{4,4}\},$$

so the required right side is **BB**, thus

$$X_{3,4} = \{A\}.$$

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|-----|-----|-----|----------|----------|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| {A} | {A} | {B} | {B} | {B} |
| a | a | b | b | b |

Likewise for $X_{4,5}$, finishing the second row.

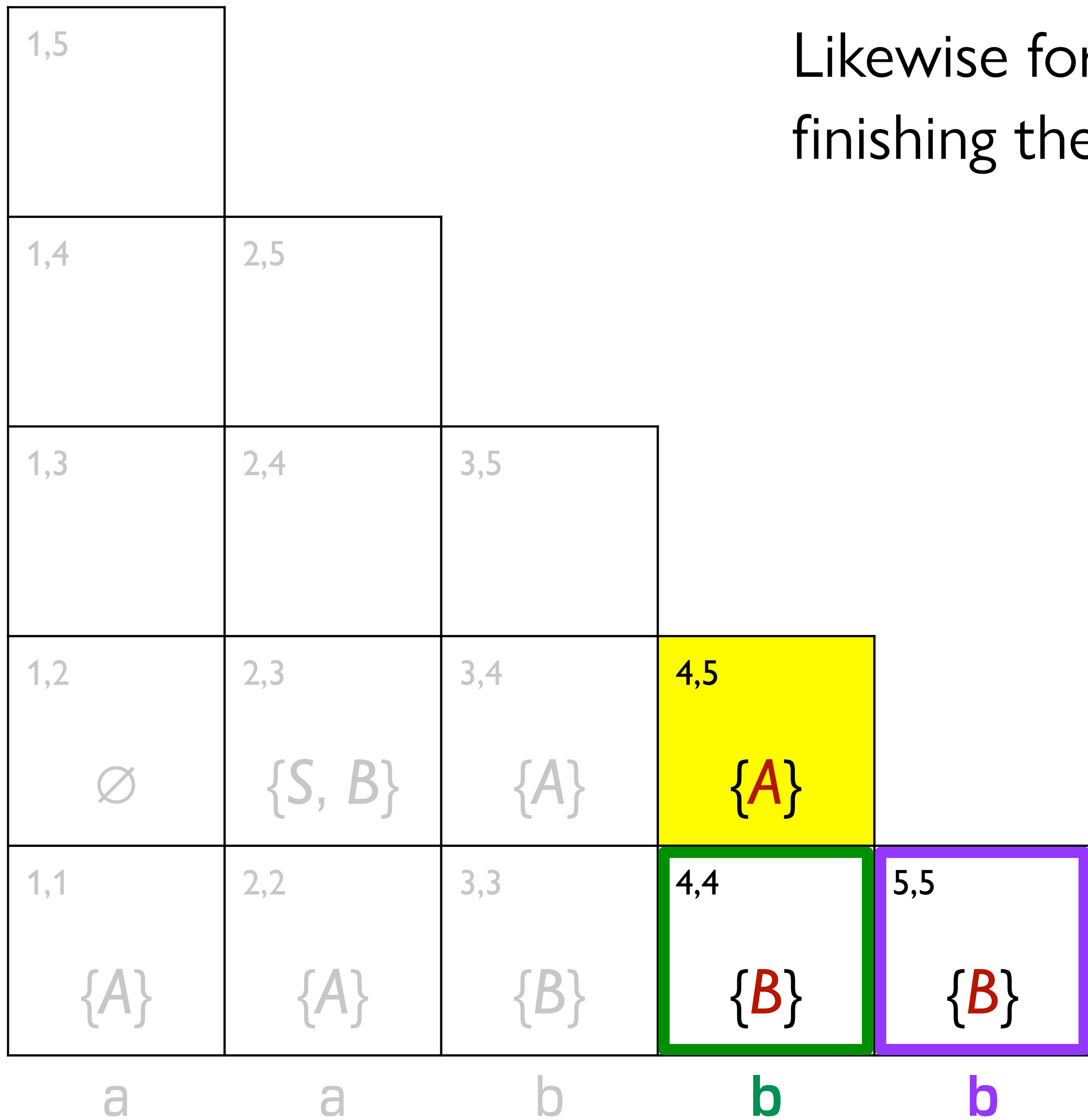
$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

Likewise for $X_{4,5}$,
finishing the second row.

| | | | | |
|-----|-----|-----|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| | | | | |
| a | a | b | b | b |

The table contains sets of strings in its cells. The cell (1,2) contains the empty set \emptyset . The cell (2,3) contains the set $\{A\}$. The cell (3,4) is highlighted in yellow. The cell (4,4) contains the set $\{B\}$ and is highlighted with a green border. The cell (5,5) contains the set $\{B\}$ and is highlighted with a purple border.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



Likewise for $X_{4,5}$, finishing the second row.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|-------------|------------|---------|---------|---------|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| \emptyset | $\{S, B\}$ | $\{A\}$ | $\{A\}$ | |
| $\{A\}$ | $\{A\}$ | $\{B\}$ | $\{B\}$ | $\{B\}$ |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|-------------|------------|---------|---------|---------|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| | \emptyset | $\{S, B\}$ | $\{A\}$ | $\{A\}$ | |
| | $\{A\}$ | $\{A\}$ | $\{B\}$ | $\{B\}$ | $\{B\}$ |
| | a | a | b | b | b |

For $X_{1,3}$, we need to consider all ways to cover those three characters with two non-terminals.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|-------------|-----|-----|-----|-----|
| 1,5 | | | | | |
| 1,4 | | 2,5 | | | |
| 1,3 | | 2,4 | 3,5 | | |
| 1,2 | \emptyset | 2,3 | 3,4 | 4,5 | |
| 1,1 | $\{A\}$ | 2,2 | 3,3 | 4,4 | 5,5 |
| | a | a | b | b | b |

For $X_{1,3}$, we need to consider all ways to cover those three characters with two non-terminals.

Body is AS or AB .

No AS , but S and B have AB as a body.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------------------|-------------------|----------------|----------------|----------------|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 $\{S, B\}$ | 2,4 | 3,5 | | |
| 1,2 \emptyset | 2,3 $\{S, B\}$ | 3,4 $\{A\}$ | 4,5 $\{A\}$ | |
| 1,1 $\{A\}$ | 2,2 $\{A\}$ | 3,3 $\{B\}$ | 4,4 $\{B\}$ | 5,5 $\{B\}$ |
| a | a | b | b | b |

For $X_{1,3}$, we need to consider all ways to cover those three characters with two non-terminals.

Body is AS or AB .

No AS , but S and B have AB as a body.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|----------------------|----------------------|-------------------|------------|------------|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 {S, B} | 2,4 | 3,5 | | |
| 1,2 \emptyset | 2,3 {S, B} | 3,4 {A} | 4,5 {A} | |
| 1,1 {A} | 2,2 {A} | 3,3 {B} | 4,4 {B} | 5,5 {B} |
| a | a | b | b | b |

For $X_{1,3}$, we need to consider all ways to cover those three characters with two non-terminals.

There's no non-terminal spanning 1–2, so we can't use this combination.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|----------|----------|----------|-----|---|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| | a | a | b | b | b |

Cell (1,3) contains $\{S, B\}$
 Cell (2,3) contains $\{S, B\}$
 Cell (1,2) contains \emptyset
 Cell (3,4) contains $\{A\}$
 Cell (4,5) contains $\{A\}$
 Cell (1,1) contains $\{A\}$
 Cell (2,2) contains $\{A\}$
 Cell (3,3) contains $\{B\}$
 Cell (4,4) contains $\{B\}$
 Cell (5,5) contains $\{B\}$

For $X_{1,3}$, we need to consider all ways to cover those three characters with two non-terminals.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|--------|--------|-----|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| {S, B} | | | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| ∅ | {S, B} | {A} | {A} | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|-------------|--------|-----|-----|-----|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| | a | a | b | b | b |
| | {A} | {A} | {B} | {B} | {B} |
| | \emptyset | {S, B} | {A} | {A} | |
| | {S, B} | | | | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|--------|--------|-----|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| {S, B} | | | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| ∅ | {S, B} | {A} | {A} | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------|--------|-----|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| {S, B} | | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| ∅ | {S, B} | {A} | {A} | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| {A} | {A} | {B} | {B} | {B} |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------|--------|-----|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| {S, B} | {A} | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| ∅ | {S, B} | {A} | {A} | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| {A} | {A} | {B} | {B} | {B} |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|--------|--------|-----|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| {S, B} | {A} | | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| ∅ | {S, B} | {A} | {A} | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|--------|--------|-----|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| {S, B} | {A} | | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| ∅ | {S, B} | {A} | {A} | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|-----|-----|-----|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| a | a | b | b | b |

Detailed description of the table: The table is a triangular grid of cells. The top row has one cell labeled '1,5'. The second row has two cells: '1,4' and '2,5'. The third row has three cells: '1,3', '2,4', and '3,5'. The fourth row has four cells: '1,2', '2,3', '3,4', and '4,5'. The fifth row has five cells: '1,1', '2,2', '3,3', '4,4', and '5,5'. Below each column, there is a label: 'a', 'a', 'b', 'b', 'b'. The content of each cell is as follows: (1,1) {A}, (1,2) ∅, (1,3) {S, B}, (1,4) {S, B}, (1,5) {A}; (2,1) {A}, (2,2) {A}, (2,3) {S, B}, (2,4) {A}, (2,5) {A}; (3,1) {A}, (3,2) {A}, (3,3) {B}, (3,4) {B}, (3,5) {B}; (4,1) {A}, (4,2) {A}, (4,3) {B}, (4,4) {B}, (4,5) {B}; (5,1) {A}, (5,2) {A}, (5,3) {B}, (5,4) {B}, (5,5) {B}. The cell (3,5) is highlighted in yellow. The cell (3,4) is highlighted with a green border. The cell (5,5) is highlighted with a purple border.

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------|--------|--------|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| {S, B} | {A} | {S, B} | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| ∅ | {S, B} | {A} | {A} | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| {A} | {A} | {B} | {B} | {B} |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|--------|--------|--------|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| {S, B} | {A} | {S, B} | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| ∅ | {S, B} | {A} | {A} | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

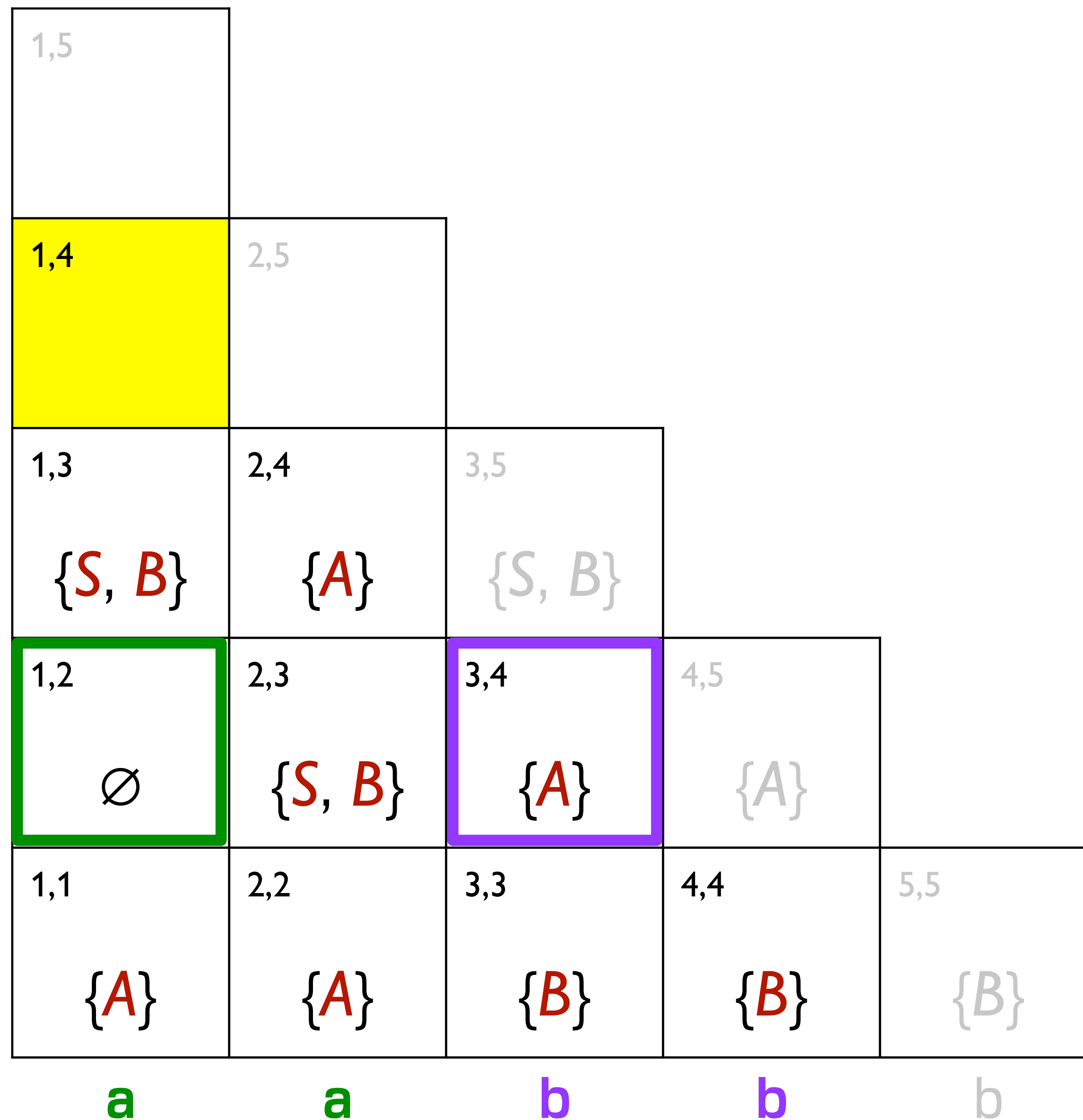
$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------|--------|--------|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| {S, B} | {A} | {S, B} | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| ∅ | {S, B} | {A} | {A} | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| {A} | {A} | {B} | {B} | {B} |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-------------|--------|--------|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| {S, B} | {A} | {S, B} | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| \emptyset | {S, B} | {A} | {A} | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$



$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|-------------|------------|------------|---------|---------|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| $\{S, B\}$ | $\{A\}$ | $\{S, B\}$ | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| \emptyset | $\{S, B\}$ | $\{A\}$ | $\{A\}$ | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| $\{A\}$ | $\{A\}$ | $\{B\}$ | $\{B\}$ | $\{B\}$ |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|----------------------|----------------------|----------------------|-------------------|-------------------|
| 1,5 | | | | |
| 1,4 {A} | 2,5 | | | |
| 1,3 {S, B} | 2,4 {A} | 3,5 {S, B} | | |
| 1,2 \emptyset | 2,3 {S, B} | 3,4 {A} | 4,5 {A} | |
| 1,1 {A} | 2,2 {A} | 3,3 {B} | 4,4 {B} | 5,5 {B} |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|----------------------|----------------------|----------------------|-------------------|-------------------|
| 1,5 | | | | |
| 1,4 {A} | 2,5 | | | |
| 1,3 {S, B} | 2,4 {A} | 3,5 {S, B} | | |
| 1,2 \emptyset | 2,3 {S, B} | 3,4 {A} | 4,5 {A} | |
| 1,1 {A} | 2,2 {A} | 3,3 {B} | 4,4 {B} | 5,5 {B} |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-------------|--------|--------|-----|-----|--|
| 1,5 | | | | | |
| 1,4 | 2,5 | | | | |
| {A} | | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| {S, B} | {A} | {S, B} | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| \emptyset | {S, B} | {A} | {A} | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {A} | {A} | {B} | {B} | {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|--------|--------|--------|-----|-----|
| 1,5 | | | | | |
| 1,4 | {A} | 2,5 | | | |
| 1,3 | {S, B} | 2,4 | 3,5 | | |
| | | {A} | {S, B} | | |
| 1,2 | ∅ | 2,3 | 3,4 | 4,5 | |
| | | {S, B} | {A} | {A} | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| | {A} | {A} | {B} | {B} | {B} |
| | a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|--------|-----|-----|-----|-----|
| 1,5 | | | | | |
| 1,4 | {A} | 2,5 | | | |
| 1,3 | {S, B} | 2,4 | 3,5 | | |
| 1,2 | ∅ | 2,3 | 3,4 | 4,5 | |
| 1,1 | {A} | 2,2 | 3,3 | 4,4 | 5,5 |
| | a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|--------|--------|--------|-----|-----|
| 1,5 | | | | | |
| 1,4 | {A} | {S, B} | | | |
| 1,3 | {S, B} | {A} | {S, B} | | |
| 1,2 | ∅ | {S, B} | {A} | {A} | |
| 1,1 | {A} | {A} | {B} | {B} | {B} |
| | a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|--------|---------------|----------|------------|----------|
| 1,5 | | | | | |
| 1,4 | {A} | {S, B} | | | |
| 1,3 | {S, B} | {A} | {S, B} | | |
| 1,2 | ∅ | {S, B} | {A} | {A} | |
| 1,1 | {A} | {A} | {B} | {B} | {B} |
| | a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|-----|--------|--------|--------|-----|
| 1,5 | | | | |
| 1,4 | {A} | {S, B} | | |
| 1,3 | {S, B} | {A} | {S, B} | |
| 1,2 | ∅ | {S, B} | {A} | {A} |
| 1,1 | {A} | {A} | {B} | {B} |
| | a | a | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|-----|--------|---------------|---------------|------------|------------|
| 1,5 | | | | | |
| 1,4 | {A} | {S, B} | | | |
| 1,3 | {S, B} | {A} | {S, B} | | |
| 1,2 | ∅ | {S, B} | {A} | {A} | |
| 1,1 | {A} | {A} | {B} | {B} | {B} |
| | a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------|--------|--------|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| {A} | {S, B} | | | |
| 1,3 | 2,4 | 3,5 | | |
| {S, B} | {A} | {S, B} | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| ∅ | {S, B} | {A} | {A} | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |
| {A} | {A} | {B} | {B} | {B} |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------------------|---------------|---------------|------------|------------|
| 1,5 | | | | |
| 1,4 {A} | 2,5 {S, B} | | | |
| 1,3 {S, B} | 2,4 {A} | 3,5 {S, B} | | |
| 1,2 \emptyset | 2,3 {S, B} | 3,4 {A} | 4,5 {A} | |
| 1,1 {A} | 2,2 {A} | 3,3 {B} | 4,4 {B} | 5,5 {B} |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|--------------------|---------------|---------------|------------|------------|--|
| 1,5 | | | | | |
| 1,4 {A} | 2,5 {S, B} | | | | |
| 1,3 {S, B} | 2,4 {A} | 3,5 {S, B} | | | |
| 1,2 \emptyset | 2,3 {S, B} | 3,4 {A} | 4,5 {A} | | |
| 1,1 {A} | 2,2 {A} | 3,3 {B} | 4,4 {B} | 5,5 {B} | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|--------------------|-------------------|-------------------|----------------|----------------|--|
| 1,5 $\{S, B\}$ | | | | | |
| 1,4 $\{A\}$ | 2,5 $\{S, B\}$ | | | | |
| 1,3 $\{S, B\}$ | 2,4 $\{A\}$ | 3,5 $\{S, B\}$ | | | |
| 1,2 \emptyset | 2,3 $\{S, B\}$ | 3,4 $\{A\}$ | 4,5 $\{A\}$ | | |
| 1,1 $\{A\}$ | 2,2 $\{A\}$ | 3,3 $\{B\}$ | 4,4 $\{B\}$ | 5,5 $\{B\}$ | |
| a | a | b | b | b | |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------------------|-------------------|-------------------|----------------|----------------|
| 1,5 $\{S, B\}$ | | | | |
| 1,4 $\{A\}$ | 2,5 $\{S, B\}$ | | | |
| 1,3 $\{S, B\}$ | 2,4 $\{A\}$ | 3,5 $\{S, B\}$ | | |
| 1,2 \emptyset | 2,3 $\{S, B\}$ | 3,4 $\{A\}$ | 4,5 $\{A\}$ | |
| 1,1 $\{A\}$ | 2,2 $\{A\}$ | 3,3 $\{B\}$ | 4,4 $\{B\}$ | 5,5 $\{B\}$ |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------------------|-------------------|-------------------|----------------|----------------|
| 1,5 $\{S, B\}$ | | | | |
| 1,4 $\{A\}$ | 2,5 $\{S, B\}$ | | | |
| 1,3 $\{S, B\}$ | 2,4 $\{A\}$ | 3,5 $\{S, B\}$ | | |
| 1,2 \emptyset | 2,3 $\{S, B\}$ | 3,4 $\{A\}$ | 4,5 $\{A\}$ | |
| 1,1 $\{A\}$ | 2,2 $\{A\}$ | 3,3 $\{B\}$ | 4,4 $\{B\}$ | 5,5 $\{B\}$ |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------------------|-------------------|-------------------|----------------|----------------|
| 1,5 $\{S, B\}$ | | | | |
| 1,4 $\{A\}$ | 2,5 $\{S, B\}$ | | | |
| 1,3 $\{S, B\}$ | 2,4 $\{A\}$ | 3,5 $\{S, B\}$ | | |
| 1,2 \emptyset | 2,3 $\{S, B\}$ | 3,4 $\{A\}$ | 4,5 $\{A\}$ | |
| 1,1 $\{A\}$ | 2,2 $\{A\}$ | 3,3 $\{B\}$ | 4,4 $\{B\}$ | 5,5 $\{B\}$ |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | |
|--------------------|-------------------|-------------------|----------------|----------------|
| 1,5 $\{S, B\}$ | | | | |
| 1,4 $\{A\}$ | 2,5 $\{S, B\}$ | | | |
| 1,3 $\{S, B\}$ | 2,4 $\{A\}$ | 3,5 $\{S, B\}$ | | |
| 1,2 \emptyset | 2,3 $\{S, B\}$ | 3,4 $\{A\}$ | 4,5 $\{A\}$ | |
| 1,1 $\{A\}$ | 2,2 $\{A\}$ | 3,3 $\{B\}$ | 4,4 $\{B\}$ | 5,5 $\{B\}$ |
| a | a | b | b | b |

$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

| | | | | | |
|--------------------|-------------------|-------------------|----------------|----------------|--|
| 1,5 $\{S, B\}$ | | | | | |
| 1,4 $\{A\}$ | 2,5 $\{S, B\}$ | | | | |
| 1,3 $\{S, B\}$ | 2,4 $\{A\}$ | 3,5 $\{S, B\}$ | | | |
| 1,2 \emptyset | 2,3 $\{S, B\}$ | 3,4 $\{A\}$ | 4,5 $\{A\}$ | | |
| 1,1 $\{A\}$ | 2,2 $\{A\}$ | 3,3 $\{B\}$ | 4,4 $\{B\}$ | 5,5 $\{B\}$ | |
| a | a | b | b | b | |

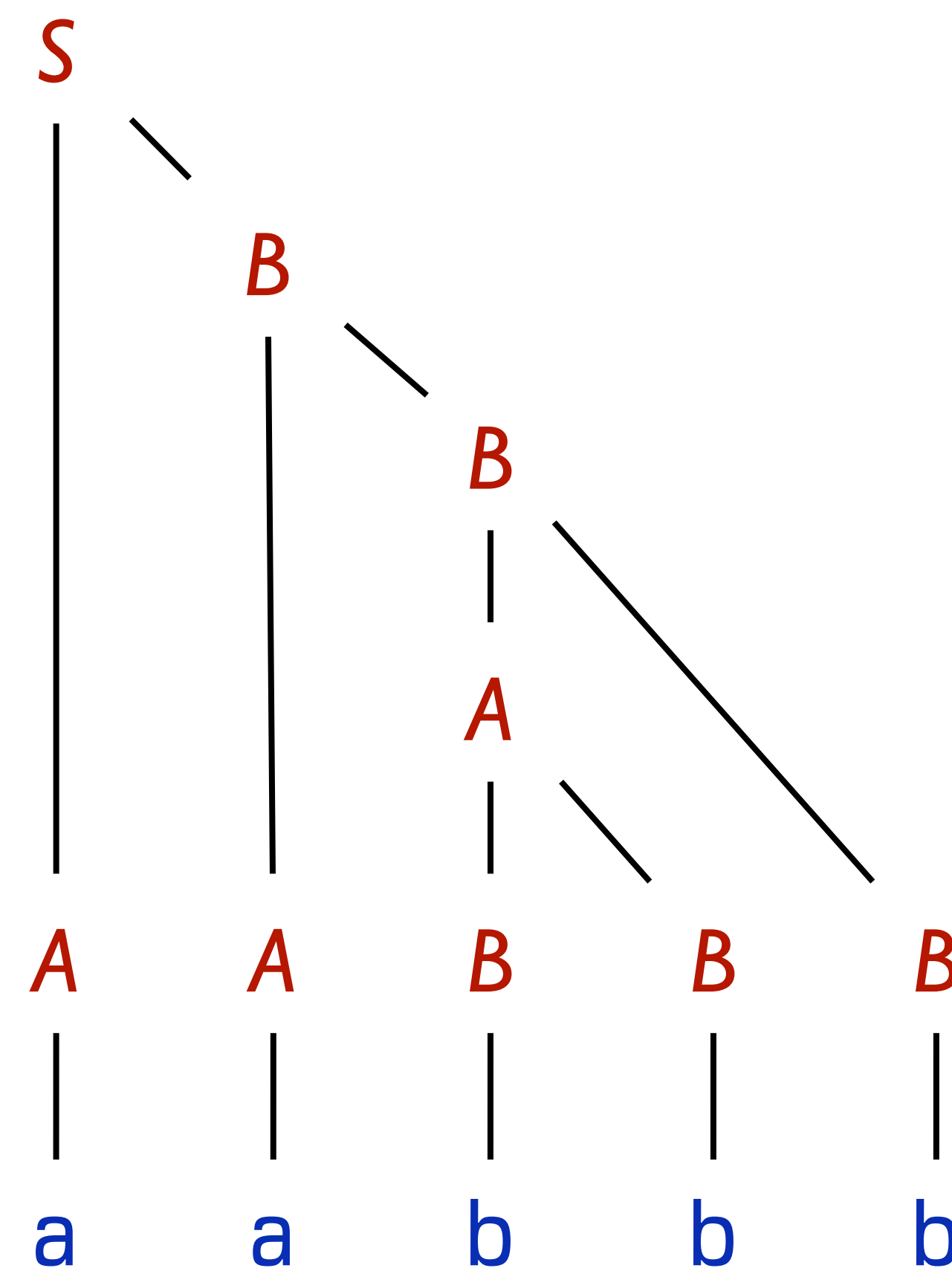
$S \rightarrow AB$
 $A \rightarrow BB \mid a$
 $B \rightarrow AB \mid b$

Since $S \in X_{1,5}$, $w \in L(G)$

| | | | | |
|--------------------|-------------------|-------------------|----------------|----------------|
| 1,5 $\{S, B\}$ | | | | |
| 1,4 $\{A\}$ | 2,5 $\{S, B\}$ | | | |
| 1,3 $\{S, B\}$ | 2,4 $\{A\}$ | 3,5 $\{S, B\}$ | | |
| 1,2 \emptyset | 2,3 $\{S, B\}$ | 3,4 $\{A\}$ | 4,5 $\{A\}$ | |
| 1,1 $\{A\}$ | 2,2 $\{A\}$ | 3,3 $\{B\}$ | 4,4 $\{B\}$ | 5,5 $\{B\}$ |
| a | a | b | b | b |

The parse tree for the input is the selection of rules that let us cover the entire string with the start symbol.

| | | | | |
|--------------------|---------------|---------------|------------|------------|
| 1,5 {S, B} | | | | |
| 1,4 {A} | 2,5 {S, B} | | | |
| 1,3 {S, B} | 2,4 {A} | 3,5 {S, B} | | |
| 1,2 \emptyset | 2,3 {S, B} | 3,4 {A} | 4,5 {A} | |
| 1,1 {A} | 2,2 {A} | 3,3 {B} | 4,4 {B} | 5,5 {B} |
| a | a | b | b | b |



CYK as a parsing algorithm

Applicability of the CYK algorithm as a parser is limited by the computational requirements needed to find a derivation.

For an input string of length n , $(n^2+n)/2$ sets need to be constructed to complete the dynamic programming table.

Each of these sets may require the consideration of several decompositions of the associated substring.

There are other parsing algorithms we can use!

Putting it all together

Parsing techniques like CYK get used with grammars for programming languages – and human languages.

Let $\Sigma = \{a, \text{Alice}, \text{big}, \text{cat}, \text{chased}, \text{dog}, \text{saw}, \text{the}\}$.

$S \rightarrow NP VP$

$NP \rightarrow Det Adj N \mid Name$

$VP \rightarrow V NP \mid V$

$Adj \rightarrow \text{big} \mid \epsilon$

$Det \rightarrow \text{the} \mid a$

$N \rightarrow \text{cat} \mid \text{dog}$

$V \rightarrow \text{chased} \mid \text{saw}$

$Name \rightarrow \text{Alice}$

Is $\text{Alice saw the big cat}$ generated by this grammar?

Let $\Sigma = \{a, \text{Alice}, \text{big}, \text{cat}, \text{chased}, \text{dog}, \text{saw}, \text{the}\}$.

$S \rightarrow NP VP$

$NP \rightarrow Det Adj N \mid Name$

$VP \rightarrow V NP \mid V$

$Adj \rightarrow \text{big} \mid \epsilon$

$Det \rightarrow \text{the} \mid a$

$N \rightarrow \text{cat} \mid \text{dog}$

$V \rightarrow \text{chased} \mid \text{saw}$

$Name \rightarrow \text{Alice}$

Step 1: Convert to CNF

1. Make start symbol non-recursive
2. Eliminate all ϵ -rules (except $S \rightarrow \epsilon$)
3. Eliminate all chain/unit rules
4. Eliminate useless symbols
5. Split rules if necessary

Is $\text{Alice saw the big cat}$ generated by this grammar?

Let $\Sigma = \{a, \text{Alice}, \text{big}, \text{cat}, \text{chased}, \text{dog}, \text{saw}, \text{the}\}$.

$S \rightarrow NP VP$

$NP \rightarrow Det X \mid Det N \mid \text{Alice}$

$X \rightarrow Adj N$

$VP \rightarrow V NP \mid \text{chased} \mid \text{saw}$

$Adj \rightarrow \text{big}$

$Det \rightarrow \text{the} \mid a$

$N \rightarrow \text{cat} \mid \text{dog}$

$V \rightarrow \text{chased} \mid \text{saw}$

Is $\text{Alice saw the big cat}$ generated by this grammar?

| | | | | |
|-----|-----|-----|-----|-----|
| 1,5 | | | | |
| 1,4 | 2,5 | | | |
| 1,3 | 2,4 | 3,5 | | |
| 1,2 | 2,3 | 3,4 | 4,5 | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 |

Step 2: CYK parsing

Alice

saw

the

big

cat

| | | | | | |
|-------|---------|-------|-------|-----|--|
| 1,5 | | | | | |
| {S} | | | | | |
| 1,4 | 2,5 | | | | |
| ∅ | {VP} | | | | |
| 1,3 | 2,4 | 3,5 | | | |
| ∅ | ∅ | {NP} | | | |
| 1,2 | 2,3 | 3,4 | 4,5 | | |
| {S} | ∅ | ∅ | {X} | | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | |
| {NP} | {V, VP} | {Det} | {Adj} | {N} | |
| Alice | saw | the | big | cat | |

Let $\Sigma = \{\text{if}, \text{elif}, \text{else:}, :, \text{True}, \text{False}\}$.

Expr \rightarrow *if Bool : Expr Mid End* | *Bool*

Mid \rightarrow *elif Bool : Expr Mid* | ϵ

End \rightarrow *else: Expr* | ϵ

Bool \rightarrow *True* | *False*

Is *if True : False else: True* generated by this grammar?

Let $\Sigma = \{\text{if}, \text{elif}, \text{else:}, :, \text{True}, \text{False}\}$.

Expr \rightarrow **if** *Bool* : *Expr* *Mid* *End* | *Bool*

Mid \rightarrow elif *Bool* : *Expr* *Mid* | ϵ

End \rightarrow else: *Expr* | ϵ

Bool \rightarrow True | False

To make this more manageable, let's remove the rules about elif.

Is **if True : False else: True** generated by this grammar?

Let $\Sigma = \{\text{if}, \text{else:}, :, \text{True}, \text{False}\}$.

$\text{Expr} \rightarrow \text{if Bool : Expr End} \mid \text{Bool}$

$\text{End} \rightarrow \text{else: Expr} \mid \epsilon$

$\text{Bool} \rightarrow \text{True} \mid \text{False}$

Step 1: Convert to CNF

1. Make start symbol non-recursive
2. Eliminate all ϵ -rules (except $S \rightarrow \epsilon$)
3. Eliminate all chain/unit rules
4. Eliminate useless symbols
5. Split rules if necessary

Is `if True : False else: True` generated by this grammar?

| | | | | | |
|-----|-----|-----|-----|-----|-----|
| 1,6 | | | | | |
| 1,5 | 2,6 | | | | |
| 1,4 | 2,5 | 3,6 | | | |
| 1,3 | 2,4 | 3,5 | 4,6 | | |
| 1,2 | 2,3 | 3,4 | 4,5 | 5,6 | |
| 1,1 | 2,2 | 3,3 | 4,4 | 5,5 | 6,6 |

Step 2: CYK parsing

if

True

:

False

else:

True

Appendix: CNF conversion details

Conversion to CNF

1. Make start symbol non-recursive
2. Eliminate all ϵ -rules (except $S \rightarrow \epsilon$)
3. Eliminate all chain/unit rules (e.g., $A \rightarrow B$)
4. Eliminate useless symbols
5. Split rules if necessary

Non-recursive start symbol

The first goal is to limit the role of the start symbol S to the initiation of a derivation.

To produce a CFG G' that is equivalent to G but does not have S in the body of any rules:

If S does not appear in the body of any rule, then $G' = G$

Otherwise, add a new start variable S' and the rule $S' \rightarrow S$

Eliminating ϵ -rules

We want to eliminate variables that don't generate terminal symbols.

Advantage: Reduces the length of derivations

Disadvantage: Increases the number of rules

This requires the identification of a set of *nullable* variables, i.e., variables that can derive the empty string by a sequence of rule applications, $A \xRightarrow{*} \epsilon$.

Eliminating ϵ -rules: Finding nullable variables

1. Put each variable A for which a rule $A \rightarrow \epsilon$ exists in the set *NULLABLE*.

2. Repeat:

For each variable $A \in V$:

If $A \rightarrow A_1A_2\dots A_k$ is a rule where each A_i is in *NULLABLE*, then add A to *NULLABLE*

If nothing was added to *NULLABLE* in this iteration, stop.

Eliminating ϵ -rules

1. Construct the set of nullable variables, as described.
2. If the start symbol S is in the nullable set, add the rule $S \rightarrow \epsilon$.
3. Consider each rule $A \rightarrow \omega$, where ω contains a sequence of n nullable variables, i.e., $\omega = \omega_1 A_1 \omega_2 A_2 \dots \omega_n A_n \omega_{n+1}$, where each ω_i is a sequence of terminals and non-nullable variables.

Replace $A \rightarrow \omega$ with a set of 2^n rules, where each new rule contains the original sequence of terminals and non-nullable variables combined with all possible combinations of the nullable variables.

4. Remove all $A \rightarrow \epsilon$ rules, except $S \rightarrow \epsilon$.

Eliminating ϵ -rules: Example

Consider the rule $A \rightarrow BABa$ and suppose B is a nullable variable.

Since there are two nullable variables in the rule body, we must replace this one rule with the original rule and three new rules (total: $2^2 = 4$ rules):

$$A \rightarrow BABa$$

$$A \rightarrow ABa$$

$$A \rightarrow BAa$$

$$A \rightarrow Aa$$

Eliminating chain rules

A rule of the form $A \rightarrow B$ is called a chain rule or a unit production.

These rules are just renaming variables in sentential forms; they mean the same strings that are derivable from B are also derivable from A .

Since chain rules do not make progress toward deriving a terminal string, we can remove them from the CFG.

More rules are added in the process, to make sure the same language is derivable.

Eliminating chain rules

Another step to reduce the length of derivations (while possibly increasing the number of rules):

For each variable A :

Compute $CHAIN(A)$, the set of all variables derivable from A by chain rules.

For each non-chain rule $B \rightarrow \omega$:

If $B \in CHAIN(A)$, replace rule $A \rightarrow B$ with $A \rightarrow \omega$

Computing $CHAIN(A)$

$$CHAIN(A) = \{A\}$$

Repeat:

For each variable B added to $CHAIN(A)$ last cycle:

For each rule $B \rightarrow C$:

Add C to $CHAIN(A)$

If $CHAIN(A)$ was not changed in this iteration, stop.

Eliminating useless symbols

A variable is *useless* if it never appears in the derivation of a string. To eliminate useless symbols:

1. *Eliminate variables that cannot derive terminal strings.* To do this, find all variables that *can* derive a terminal string.

$TERM = \{A \mid A \rightarrow \text{string of terminals is a rule}\}$

Repeat:

For each $A \in TERM$:

If $A \rightarrow \omega$ is a rule, where ω is a string of terminals and variables in $TERM$:

Add A to $TERM$

If no more variables were added to $TERM$ in this iteration, stop.

Then eliminate:

- a. variables not in $TERM$
- b. rules involving a variable not in $TERM$

Eliminating useless symbols

2. *Eliminate symbols that are not reachable from S.* To do this, find all variables that are reachable:

$$REACH = \{S\}$$

Repeat:

For each A just added to $REACH$:

For each rule $A \rightarrow \omega$:

Add all variables in ω to $REACH$

If $REACH$ was not changed, stop.

Then eliminate:

- a. variables not in $REACH$
- b. rules involving a variable not in $REACH$
- c. all terminals not in remaining rules

At this point, all rules are of the form

$$S \rightarrow \epsilon$$

$$A \rightarrow a$$

$A \rightarrow \omega$, where ω is a string of more than one terminals, variables, or terminals and variables (excluding S)

At this point, all rules are of the form

$$S \rightarrow \varepsilon$$

$$A \rightarrow a$$

$A \rightarrow \omega$, where ω is a string of more than one terminals, variables, or terminals and variables (excluding S)

But Chomsky normal form is even stricter, allowing only:

$$S \rightarrow \varepsilon$$

$$A \rightarrow a$$

$A \rightarrow BC$, where neither $B \neq S$ and $C \neq S$

Split rules

Step 1: For each rule $A \rightarrow \omega$, replace each terminal a in ω with a new variable V_a . Add the rule $V_a \rightarrow a$.

E.g.,

$$A \rightarrow bDcF$$

becomes

$$A \rightarrow V_b D V_c F$$

$$V_b \rightarrow b$$

$$V_c \rightarrow c$$

Split rules

Step 2: Now introduce new variables to remove variables one by one to get only two variables in the body of each rule.

E.g.,

$$A \rightarrow BCDF$$

becomes

$$A \rightarrow BT_1$$

$$T_1 \rightarrow CT_2$$

$$T_2 \rightarrow DF$$

