# Towards Relational Contextual Equality Saturation

TYLER HOU, UC Berkeley, USA
SHADAJ LADDAD, UC Berkeley, USA
JOSEPH M. HELLERSTEIN, UC Berkeley, USA

Equality saturation is a powerful technique for program optimization. Contextual equality saturation extends this to support rewrite rules that are conditioned on where a term appears in an expression. Existing work has brought contextual reasoning to egg; in this paper, we share our ongoing work to extend this to relational equality saturation in egglog. We summarize the existing approaches to contextual equality saturation, outline its main applications, and identify key challenges in combining this approach with relational models.

Additional Key Words and Phrases: e-graphs, program analysis, query optimization

## 1 INTRODUCTION

In recent years there as been a surge of interest in using equality saturation to optimize programs [Tate et al. 2009]. Equality saturation is an optimization technique which repeatedly applies rewrite rules to an program. These rewrite rules add equalities between (sub)terms of the expression to a data structure called an e-graph. The e-graph represents a possibly-infinite set programs equivalent to the original program. Finally, an optimized program can be extracted from the e-graph.

However, existing equality saturation tools like egg [Willsey et al. 2021] and egglog [Zhang et al. 2023] do not natively support *contextual equalities*. Contextual equalities are equalities between terms that are valid in some subgraph, but may not be valid in general. For example, consider the ternary $x == 2 ? (x \times y) : y$. Under the then branch of the ternary, the condition $x == 2$ is satisfied, so an optimizer could replace the multiplication with a bitshift, producing an optimized expression $x == 2 ? (y \gg 1) : y$. However, equating the term $x \times y$ with $y \gg 1$ in general is not valid, because $x \times y$ may appear in some other term where the condition is not satisfied.

Reasoning under contexts is necessary in many program optimizers to achieve better performance. For example, by applying context-aware optimizations to circuit design, an RTL optimization tool reduced circuit area by 41% and delay by 33% [Coward et al. 2023]. In addition, reasoning about contextual properties like sort order is critical for dataflow optimization and relational query optimization [Graefe 1995; Laddad et al. 2023].

## 2 EXISTING APPROACHES

There have been various approaches to support contextual reasoning in equality saturation frameworks. The authors of the RTL optimization tool encoded contextual reasoning in egg by adding ASSUME(x, c) e-nodes to the e-graph, where x is the expression to be optimized, and c is a set of expressions which represent known constraints [Coward et al. 2023]. Auxiliary rewrite rules "push down" ASSUME e-nodes into the expression tree, adding additional constraints when appropriate. Finally, after ASSUME e-nodes have been pushed down sufficiently far, domain-specific rules rewrite them into more efficient terms using the constraints. However, the main limitation of this approach is that these "extra" ASSUME e-nodes significantly expand the size of the e-graph, harming performance. In their tool, optimization of a floating point subtactor took 22 minutes, with the majority of time spent in e-graph expansion [Coward et al. 2023].

An alternative to adding explicit ASSUME e-nodes into the e-graph is to annotate e-classes with the contexts that they appear in via a top-down analysis. Then, a contextual rewrite rule can "copy" the contextual subgraph, replacing subterms that can be contextually rewritten with their equivalents [Drewery 2022]. However, this approach can also substantially expand the e-graph: to avoid false equalities, such a rewrite must make independent copies of all contextually-rewritten e-classes and their ancestors, up to and including the "source" of the context. When multiple contexts are nested, this can again lead to a combinatorial explosion of the e-graph size. Finally, copying the e-graph also makes it more expensive to apply context insensitive rewrites – such rewrites may have to match all copied subgraphs since they are no longer related to the original subgraph.

A promising approach for efficiently supporting contextual equalities are colored e-graphs [Singher and Itzhaky 2023]. Instead of adding new e-nodes to the e-graph, colored e-graphs support multiple equivalence relations: a base equivalence relation represents equalities that apply in every context, and context-sensitive, colored equivalences are *layered* on top of the base relation. These layered relations can be thought of "shallow" copies of the base relation. This approach has two advantages: first, it saves memory because context-sensitive rewrites typically add relatively few equivalences on top of the base equivalence. Second, when additional base equivalences are found, the context-sensitive equivalence relations can be efficiently updated. However, one limitation of the colored e-graph work is that it targets the egg library, which implements non-relational e-matching.

## 3   CASE STUDIES

To better understand the space of applications that can benefit from contextual equality saturation, we explore three case studies from a range of optimization domains.

### 3.1   Relational query optimization

Database engines work by translating user-provided queries (often in SQL) down to a query plan, then repeatedly applying rewrite rules to find more efficient plans. For example, consider the SQL query in Fig. 1, which finds all nodes reachable in a path of exactly three steps from a source node with ID 100. The initial plan uses merge joins, which is reasonable if we assume that the relations are sorted or indexed on the join columns. But there is room for improvement, since the selection $\sigma_{l.source=100}$ can be pushed down past the joins. The left side of Fig. 2 shows a plan where $\sigma_{l.source=100}$ has been pushed down as far as possible. This optimization can be discovered through equality saturation with a rewrite rule that appropriately commutes selections with joins.

```
SELECT  r.target
FROM      edges AS l
    JOIN edges AS m ON l.target = m.source
    JOIN edges AS r ON m.target = r.source
WHERE l.source = 100;
```
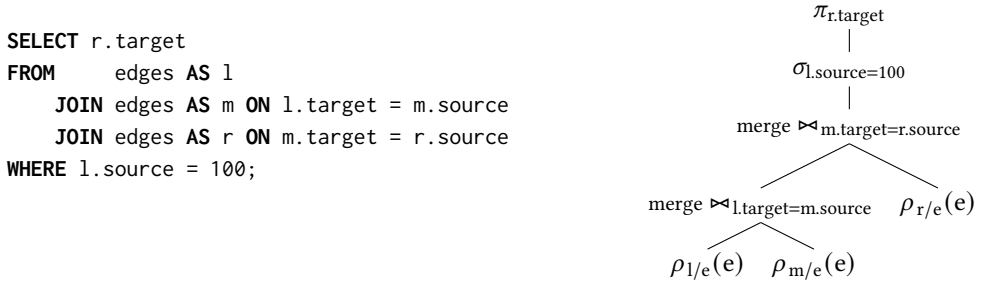


Fig. 1. An example SQL query searching for paths of length three and its corresponding initial query plan.

However, such a plan may *still* not be optimal; after selection, the inputs to the bottom join may be relatively small. If relations are small and fit in memory, then it may be more efficient to hash
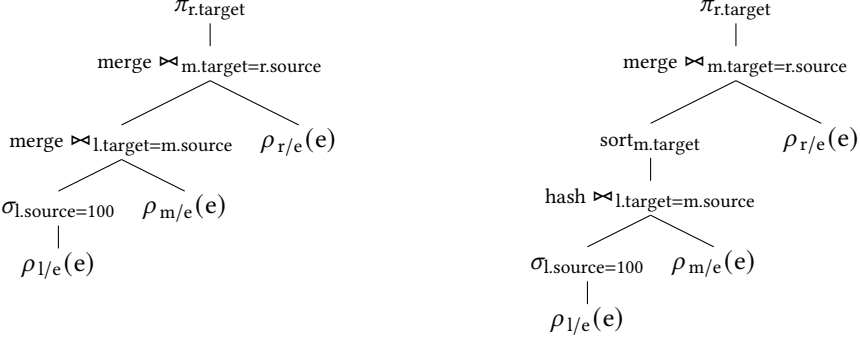
Fig. 2. Optimized versions of the default query plan in Fig. 1. In the left plan, the selection on l.source = 100 has been pushed down past the joins. In the right plan, the bottom merge join on l and m has been replaced with a hash join; a sort has been added above to enforce that the inputs to the top merge join are sorted.

join them. Thus, we would like to replace the merge join merge $\bowtie_{\text{l.target=m.source}}$ with a hash join hash $\bowtie_{\text{l.target=m.source}}$. But unconditionally rewriting merge joins to hash joins is not valid—merge joins preserve sort order, but hash joins do not. For example, if a merge join $J$ is being fed to an operator which requires that its inputs be sorted (e.g. another merge join, as in the left plan in Fig. 2, or a non-commutative aggregation), it is not valid to replace $J$ with a hash join.

It is not difficult to fix this though: we can add a sort *enforcer* [Graefe and McKenna 1993] underneath the top merge join that sorts its left input on the m.target column. Underneath the sort enforcer, it is *contextually valid* to rewrite merge joins to hash joins, because all sub-plans will be eventually sorted on the m.target column. The right hand query plan in Fig. 2 shows the query plan with the sort enforcer added and bottom merge join replaced with a hash join.

More generally, in database query optimization, query (sub-)plans must often satisfy certain *physical properties*, because their outputs are fed to consumers who rely on those physical properties [Graefe and McKenna 1993]. Physical properties include sort order, partitioning, and data location. In general, valid rewrites must preserve physical properties. However, in certain contexts, like underneath a sort enforcer, additional rewrites that "break" physical properties become valid, because those physical properties will be re-established by the enforcer.

## 3.2 Simplifying conditionals

As shown in the introduction, contextual equality saturation can also be used to discover rewrites for conditional statements. The ternary $x == 2$ ? $(x \times y)$ : $y$ can be rewritten to the ternary $x == 2$ ? $(y \gg 1)$ : $y$, since under the if branch of the ternary, we know that $x == 2$, which yields the chain of equalities $x \times y \equiv 2 \times y \equiv y \gg 2$.

But not only can contextual equality saturation reason under conditionals, in some circumstances, it can remove conditional statements entirely. Consider another ternary $(a > b)$ ? $(a > b)$ : $(a \leq b)$. With contextual equivalences, one might reason:

(1) Under the then branch of the ternary, we have the additional equivalence $a > b \equiv$ true. Hence, the e-class inside the then branch contains the terms $\{a > b, \text{true}\}$.

(2) Under the then branch of the ternary, we have the additional equivalence $\neg(a > b) \equiv$ true. Another (general) rewrite rule could rewrite $\neg(a > b) \equiv b \leq a$, so, by transitivity, $b \leq a \equiv$ true. The e-class inside the else branch would contain the terms $\{\neg(a > b), \text{true}, b \leq a\}$.

(3) An "intersecting" rule says if the bodies of both branches is equivalent to some term $t$, then the entire ternary is equivalent to $t$. That is, because *both the e-classes* for the then branch and the else branch contain the term true, then the entire ternary is equivalent to true: $\big((a > b) \ ? \ (a > b) : (a \le b)\big) \equiv$ true. (This rule is corresponds to the "disjunction elimination" or "proof by cases" inference rule in propositional logic [Singher and Itzhaky 2023].)

We note that step 3 "converts" contextual equivalences within the branches of the ternary into context-insensitive equivalences by taking the "intersection" of two equivalence relations. In order to support such reasoning, a contextual equality saturation framework needs to let a user write rules that can refer to and manipulate (contextual) equivalence relations.

## 3.3 Lambda application

Finally, we show an example where contextual reasoning can simplify the body of a lambda application. This type of contextual reasoning is similar to conditional simplification, except has one additional complication. Consider the lambda application $(\lambda x.x + 1)2$. We want to show that the entire term is equal to 3. Because 2 is applied to $\lambda x.x + 1$, we can contextually equate $x \equiv 2$ within the body of the lambda. Then, using a context-insensitive rewrite rule $2 + 1 \equiv 3$, we can transitively find the equivalence $x + 1 \equiv 3$, inside the lambda body. This means that the entire lambda application $(\lambda x.x + 1)2 \equiv (\lambda x.3)2$.

Here is the complication: we want to simplify the entire lambda application to 3, but so far we have only simplified the body. Just like in conditional rewriting, how might we "convert" contextual equivalences on the body to context-insensitive equivalences on the entire application? One might suggest the following (incorrect) rule: if the lambda body is equal to a term $t$, then the entire lambda application is equivalent to the term $t$. This is not correct because the e-class of the body contains the term $t = x + 1$, and we *cannot* equate the whole lambda application to $x + 1$, since $x$ is "unmentionable" outside of the scope of the lambda.

The rule that we want is: if the lambda body is equal to a term $t$, and that term $t$ does not transitively mention $x$, then the entire lambda application is equal to $t$. Concretely, the e-class of the lambda body would contain the terms $\{x + 1, 2 + 1, 3\}$ when saturated. Under the above rule, we would (correctly) equate the entire lambda application to the terms in $\{2 + 1, 3\}$.

However, this type of rule is hard to implement in existing equality saturation frameworks because existing equality saturations frameworks like egg and egglog *canonicalize* the e-graph—the body of the lambda application is not a set of terms, but instead an e-class. Furthermore, there is no way to "cleave apart" an e-class to recover the terms that do not contain $x$. Further research is needed to discover how contextual equality saturation might support such an operation.[1]

## 4  A SET-THEORETIC PERSPECTIVE

In this section we present ongoing work on a set-theoretic perspective on contextual equality saturation. We aim to describe contextual equality saturation in an implementation-independent manner and to suggest how a relational contextual equality saturation framework might practically support contexts. We assume familiarity with lattices and equivalence relations; definitions for lattices and equivalence relations can be found in the Definitions. We begin with a definition of *quotient sets*, and we relate them to e-classes and a hierarchy of equivalence relations.

*Definition 4.1.* Let $A$ be a set, and let $\sim$ be an equivalence relation on $A$. The *quotient set* $A/\sim$ is a set of equivalence classes where two elements $a, b \in A$ are in the same equivalence class iff $a \sim b$.

---

[1]We note that existing equality saturations can implement beta-reduction (as well as entire interpreters for a lambda calculus) by using e-graph *analyses* to keep track of free and bound variables in an expression.

*Example 4.2.* Let $\Sigma$ be a set of function symbols, and let $R = \{f(a, b, \ldots) \mid f, a, b \ldots \in \Sigma\}$ be a set of terms. For two terms $a$ and $b$, let $a \sim b$ iff $a$ and $b$ are equivalent terms. Then $R/\sim$ is exactly the set of e-classes of $R$ under $\sim$, and the quotient map $\pi : R \to R/\sim$ sends a term to its e-class (it is exactly the `find` operation in egglog [Zhang et al. 2023]).

We note that $R/\sim$ is still not the "e-graph" stored by egglog because terms within an equivalence class in $R/\sim$ still refer to terms in $R$; that is, $R/\sim$ is not canonicalized, so it may be exponentially large. We define the canonicalized database $(R/\sim)^*$ as $(R/\sim)^* = \{\{f(\pi(a), \pi(b), \ldots), \ldots\} \mid f(a, b, \ldots), \ldots \in R/\sim\}$; i.e. we replace references to subterms with references to the subterms' equivalence classes. This shrinks the size of each equivalence class because terms in $R/\sim$ that used to reference equivalent but not identical subterms are merged in $(R/\sim)^*$.

PROPOSITION 4.3. *Let $B^+$ denote the transitive closure of $B$. The set of equivalence relations on $A$ forms a lattice, where $S \sqcap T \stackrel{def}{=} S \cap T$ and $S \sqcup T \stackrel{def}{=} (S \cup T)^+$. Furthermore, $\sim_1 \leq \sim_2$ iff $\sim_1 \subseteq \sim_2$.*

Intuitively, the meet of two equivalence relations $S$ and $T$ is another equivalence relation $U$ where $a \sim b$ if $a$ is equivalent to $b$ under *both $S$ and $T$*. The join of two equivalence relations $S$ and $T$ is another equivalence relation $V$ where $a \sim_V b$ if $a$ is equivalent to $b$ by *chaining together any equivalences of $S$ or $T$*. (Transitive closure formally captures the notion of "chaining together.")

*Definition 4.4.* Let $A$ be a set, let $\sim^A$ be the set of all equivalence relations on $A$, and let $L$ be a lattice with an element called bottom, $\bot \in L$, satisfying $\forall l \in L, \bot \leq l$. A *context-annotated equivalence relation* is a mapping $\phi : L \to \sim^A$ where $\phi$ preserves order: $l_1 \leq l_2 \implies \phi(l_1) \leq \phi(l_2)$.

Definition 4.4 suggests what structure contexts should have to be suitable for labels in contextual equality saturation—they should have the property that as one learns more information, more, not fewer, equalities become available. We associate with the base context, labeled $\bot$, the fewest equivalences $\phi(\bot)$—these are the rewrites that are allowed in every context. As one traverses up the context lattice, more information is learned, and more terms are equal.

In Section 3.1, our context lattice contained two elements: $\bot$ represented a context where sort order had to be preserved, so merge join and hash join were not equivalent. Underneath a sort operator, however, we transition up the context lattice to a context $s$ where we no longer had to preserve sort order. Thus, we verify that $\phi(\bot) \leq \phi(s)$ (in fact, the equality is strict: $\phi(\bot) < \phi(s)$).

Similarly, we (secretly) used the lattice structure of equivalence relations in Section 3.2. Under each branch of the conditional, we transitioned up the lattice: in the then branch, we transitioned from $\bot$ to $t > \bot$, where $a > b \equiv$ true. In the false branch, we transitioned up to $f > \bot$, where $\neg(a > b) \equiv$ true. When we took the intersection of the terms for both branches of the ternary, we applied meet ($\sqcap$) on the contextual equivalence relations for each branch, "updating" $\phi(\bot) \leftarrow \psi(\bot)$ where $\psi(\bot) = \phi(\bot) \sqcup (\phi(t) \sqcap \phi(f))$.

PROPOSITION 4.5. *Let $A$ be a set, and let $\sim_1$ and $\sim_2$ be two equivalence relations on $A$ where $\sim_1 \leq \sim_2$. Then we have (note $1 \iff 2 \implies 3$):*

(1) *The quotient map (which sends elements to their equivalence classes) $\pi_2 : R \to R/\sim_2$ factors through the quotient map $\pi_1 : R \to R/\sim_1$: there is a map $q : R/\sim_1 \to R/\sim_2$ such that $\pi_2 = q \circ \pi_1$ See Fig. 3.*

(2) *The quotient set $A/\sim_2$ is "more coarse" than the quotient set $A/\sim_1$ in the following sense: every equivalence class $c \in A/\sim_2$ is a union of one or more equivalence classes $C \subseteq A/\sim_1$; that is, $c = \bigcup C$.*

(3) *$A/\sim_2$ has fewer equivalence classes: $|A/\sim_2| \leq |A/\sim_1|$.*
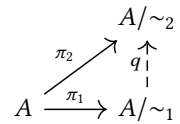


Fig. 3. $\pi_2$ factors through $\pi_1$ and $q$.

In Proposition 4.5, one can interpret $\sim_1$ as some (possibly contextual) equivalence relation and $\sim_2$ as a coarser equivalence relation. Item (1) says that any database that stores a finer database of terms $R/\sim_1$ can always "recover" the coarser set of terms by applying $q$, which sends $R/\sim_1 \to R/\sim_2$.

As observed in both egglog [Zhang et al. 2023] and colored e-graphs [Singher and Itzhaky 2023], it can be exponentially faster to e-match on a canonicalized e-graph versus an uncanonicalized e-graph. *The existence of $q$ suggests that contextual equality saturation frameworks can trade off space for time*: instead of storing a canonicalized copy of the e-graph for every equivalence relation, databases can store e-graphs for a *lower-bound* of contexts. When e-matching is requested on some context that is not explicitly materialized, frameworks have two options: either they can e-match directly on the stored e-graphs, incurring the cost of a join on the equivalence relation; or they can apply $q$ on-the-fly and e-match on a canonicalized e-graph. If $\sim_2$ adds many equivalences, then the framework will have to join on many terms, so it is likely cheaper to apply $q$ up front to canonicalize (a copy of) the e-graph to the requested equivalence. On the other hand, $\sim_2$ adds very few equivalences, then a join might be cheaper. This type of reasoning is the bread and butter of classical database management systems, which can collect statistics to estimate the size of intermediate relations and aggregations [Gray et al. 2007].

Finally, we note that item (3) in Proposition 4.5 implies that as more equivalences are added, the number of equivalence classes in the canonicalized database can shrink. Relational logic systems like Datalog rely on monotonicity to be efficient and correct. Thus, a relational equality saturation framework that supports contexts needs to prevent users from observing any non-monotonicity. For example, programs should not be able to depend on how many equivalence classes there are in the database. (egglog already achieves this for a single equivalence relation.)

## 5 CONCLUSION

We have provided an overview of existing approaches to contextual equality saturation, and shown three problems that contextual equality saturation may be able to (partially) solve. Our early work points towards a set-theoretic model for contextual equality saturation that lends itself to natural system optimizations from the database literature.

In the future, we plan to further develop our set-theoretic model and extend it to the relational setting. In addition, we aim to develop a syntax that lets users easily express context-sensitive rewrites in a relational model and manipulate equivalences across contexts, but prevents them from observing non-monotonicity (which is crucial for performance and correctness). Finally, we aim to explore whether existing Datalog systems like egglog and Soufflé can be adapted to support a hierarchy of equivalence relations.

## REFERENCES

Samuel Coward, George A. Constantinides, and Theo Drane. 2023. Automating Constraint-Aware Datapath Optimization using E-Graphs. arXiv:2303.01839 [cs.AR]

Alexandre Drewery. 2022. *Automatic Equivalence Verification for Translation Validation*. Master's thesis. ENS Rennes. https://perso.eleves.ens-rennes.fr/people/alexandre.drewery/internship_report.pdf

David S. Dummit and Richard M. Foote. 2004. *Abstract Algebra*. John Wiley & Sons, Hoboken, New Jersey.

Goetz Graefe. 1995. The Cascades Framework for Query Optimization. *IEEE Data Eng. Bull.* 18, 3 (1995), 19–29. http://sites.computer.org/debull/95SEP-CD.pdf

G. Graefe and W.J. McKenna. 1993. The Volcano optimizer generator: extensibility and efficient search. In *Proceedings of IEEE 9th International Conference on Data Engineering*. 209–218. https://doi.org/10.1109/ICDE.1993.344061

George Grätzer. 2009. *Lattice Theory: First Concepts and Distributive Lattices.* Dover Publications, Mineola, New York.

Jim Gray, Surajit Chaudhuri, Adam Bosworth, Andrew Layman, Don Reichart, Murali Venkatrao, Frank Pellow, and Hamid Pirahesh. 2007. Data Cube: A Relational Aggregation Operator Generalizing Group-By, Cross-Tab, and Sub-Totals. arXiv:cs/0701155 [cs.DB]

Shadaj Laddad, Conor Power, Tyler Hou, Alvin Cheung, and Joseph M. Hellerstein. 2023. Optimizing Stateful Dataflow with Local Rewrites. arXiv:2306.10585 [cs.PL]

Eytan Singher and Shachar Itzhaky. 2023. Colored E-Graph: Equality Reasoning with Conditions. arXiv:2305.19203 [cs.PL]

Ross Tate, Michael Stepp, Zachary Tatlock, and Sorin Lerner. 2009. Equality saturation: a new approach to optimization. In *Proceedings of the 36th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Savannah, GA, USA) *(POPL '09).* Association for Computing Machinery, New York, NY, USA, 264–276. https://doi.org/10.1145/1480881.1480915

Max Willsey, Chandrakana Nandi, Yisu Remy Wang, Oliver Flatt, Zachary Tatlock, and Pavel Panchekha. 2021. Egg: Fast and Extensible Equality Saturation. *Proc. ACM Program. Lang.* 5, POPL, Article 23 (January 2021), 29 pages. https://doi.org/10.1145/3434304

Yihong Zhang, Yisu Remy Wang, Oliver Flatt, David Cao, Philip Zucker, Eli Rosenthal, Zachary Tatlock, and Max Willsey. 2023. Better Together: Unifying Datalog and Equality Saturation. *Proc. ACM Program. Lang.* 7, PLDI, Article 125 (June 2023), 25 pages. https://doi.org/10.1145/3591239

## A DEFINITIONS

*Definition A.1.* A *meet semilattice* is an algebraic structure $(S, \sqcap)$ where $\sqcap$ (pronounced meet, or greatest lower bound) is a binary operator that is idempotent, commutative, and associative. Given a join semilattice, we can define $a \leq b$ if $a \sqcap b = a$ [Grätzer 2009].

*Definition A.2.* Dually, a *join semilattice* is an algebraic structure $(S, \sqcup)$ where $\sqcup$ (pronounced join, or least upper bound) is a binary operator that is idempotent, commutative, and associative. Given a join semilattice, we can define $a \geq b$ if $a \sqcup b = a$ [Grätzer 2009].

*Definition A.3.* A *lattice* is an algebraic structure $(S, \sqcap, \sqcup)$ where $(S, \sqcap)$ is a meet semilattice, $(S, \sqcup)$ is a join semilattice, and $\sqcap$ and $\sqcup$ satisfy the absorption law $a \sqcap (a \sqcup b) = a \sqcup (a \sqcap b) = a$.

*Example A.4.* Let $S = \mathbb{Z}$, $\sqcap = \min$, and $\sqcup = \max$. Then $(\mathbb{Z}, \min, \max)$ is a lattice. For all $a, b \in \mathbb{Z}$:
  (1) Both min and max are idempotent, associative, and commutative.
  (2) The meet of $a$ and $b$ is $\min(a, b)$, and if $\min(a, b) = a$ then $a \leq b$ (under the normal ordering).
  (3) The join of $a$ and $b$ is $\max(a, b)$, and if $\max(a, b) = a$ then $a \geq b$.
  (4) $\min(a, \max(a, b)) = a = \max(a, \min(a, b))$.

*Definition A.5.* A *binary relation* on a set $A$ is a subset $R$ of $A \times A$. For $a, b \in A$, we write $a \sim b$ if $(a, b) \in R$. A binary relation is an *equivalence relation* if it is reflexive, symmetric, and transitive [Dummit and Foote 2004].