

Affine Algebraic Decision Diagrams (AADDs) and their Application to Structured Probabilistic Inference

Scott Sanner

Department of Computer Science
University of Toronto
Toronto, ON M5S 3H5, CANADA
ssanner@cs.toronto.edu

David McAllester

TTI at Chicago
1427 East 60th Street
Chicago, IL 60637, USA
mcallester@tti-c.org

Abstract

We propose an affine extension to ADDs (AADD) capable of compactly representing context-specific, additive, and multiplicative structure. We show that the AADD has worst-case time and space performance within a multiplicative constant of that of ADDs, but that it can be linear in the number of variables in cases where ADDs are exponential in the number of variables. We provide an empirical comparison of tabular, ADD, and AADD representations used in standard Bayes net and MDP inference algorithms and conclude that the AADD performs at least as well as the other two representations, and often yields an exponential performance improvement over both when additive or multiplicative structure can be exploited. These results suggest that the AADD is likely to yield exponential time and space improvements for a variety of probabilistic inference algorithms that currently use tables or ADDs.

1 Introduction

Algebraic decision diagrams (ADDs) [1] provide an efficient means for representing and performing arithmetic operations on functions from a factored boolean domain to a real-valued range (i.e., $\mathbb{B}^n \rightarrow \mathbb{R}$). They rely on two main principles to do this:

1. ADDs represent a function $\mathbb{B}^n \rightarrow \mathbb{R}$ as a directed acyclic graph – essentially a decision tree with re-convergent branches and real-valued terminal nodes.
2. ADDs enforce a strict variable ordering on the decisions from the root to the terminal node, enabling a minimal, canonical diagram to be produced for a given function. Thus, two identical functions will always have identical ADD representations under the same variable ordering.

As shown in Figure 1, ADDs often provide an efficient representation of functions with context-specific independence [2], such as functions whose structure is conjunctive (1a) or disjunctive (1b) in nature. Thus,

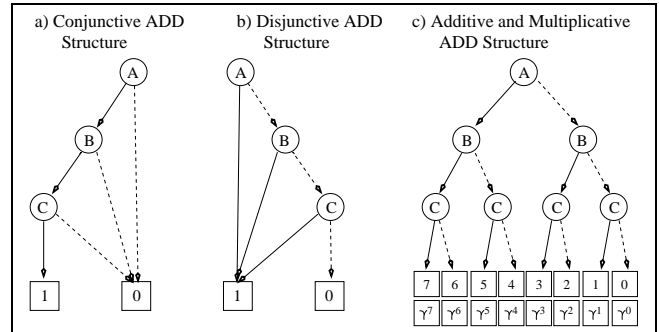


Figure 1: Some example ADDs showing a) conjunctive structure ($f = if(a \wedge b \wedge c, 1, 0)$), b) disjunctive structure ($f = if(a \vee b \vee c, 1, 0)$), and c) additive ($f = 4a + 2b + c$) and multiplicative ($f = \gamma^{4a+2b+c}$) structure (top and bottom sets of terminal values, respectively). The high edge is solid, the low edge is dotted.

ADDs can offer exponential space savings over a fully enumerated tabular representation. However, the compactness of ADDs does not extend to the case of additive or multiplicative independence, as demonstrated by the exponentially large representations when this structure is present (1c). Unfortunately such structure often occurs in probabilistic and decision-theoretic reasoning domains, leading to exponential running times and space requirements for inference in these domains.

2 Affine Algebraic Decision Diagrams

To address the limitations of ADDs, we introduce an affine extension to the ADD (AADD) that is capable of canonically and compactly representing context-specific, additive, and multiplicative structure in functions from $\mathbb{B}^n \rightarrow \mathbb{R}$.

We define AADDs with the following BNF:

$$\begin{aligned}
 F & ::= 0 \mid if(F^{var}, c_h + b_h F_h, c_l + b_l F_l) \\
 G & ::= c + bF
 \end{aligned}$$

Here we have c_h and c_l are real (or floating-point) constants in the closed interval $[0, 1]$, and b_h and b_l are real constants in the half-open interval $(0, 1]$. We also impose the following constraints:

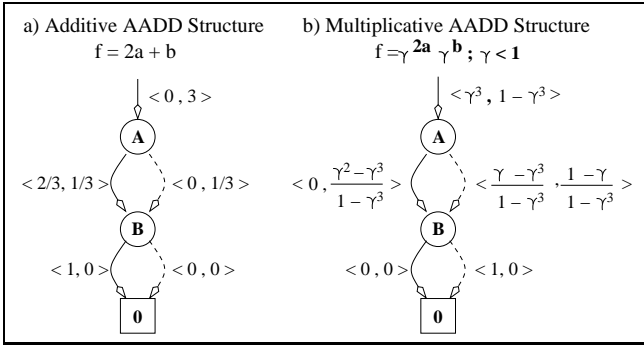


Figure 2: Portions of the ADDs from figure 1c expressed as generalized AADDs. The edge weights are given as $\langle c, b \rangle$.

1. The variable F^{var} does not appear in F_h or F_l .
2. $\min(c_h, c_l) = 0$
3. $\max(c_h + b_h, c_l + b_l) = 1$
4. If $F_h = 0$ then $b_h = 0$ and $c_h > 0$. Similarly for F_l .
5. In the grammar for G , we require that if $F = 0$ then $b = 0$, otherwise $b > 0$.

Expressions in the grammar for F will be called normalized AADDs and expressions in the grammar for G will be called generalized AADDs.¹

Let $Val(F, \rho)$ be the value of AADD F under variable value assignment ρ . This can be defined recursively by the following equation:

$$Val(F, \rho) = \begin{cases} F = 0 : & 0 \\ F \neq 0 \wedge \rho(F^{var}) = true : & c_h + b_h \cdot Val(F_h, \rho) \\ F \neq 0 \wedge \rho(F^{var}) = false : & c_l + b_l \cdot Val(F_l, \rho) \end{cases}$$

Lemma 2.1. *For any normalized AADD F we have that $Val(F, \rho)$ is in the interval $[0, 1]$, $\min_{\rho} Val(F, \rho) = 0$, and if $F \neq 0$ then $\max_{\rho} Val(F, \rho) = 1$.*

We say that F satisfies a given variable ordering if $F = 0$ or F is of the form $if(F^{var}, c_h + b_h F_h, c_l + b_l F_l)$ where F^{var} does not occur in F_h or F_l and F^{var} is the earliest variable under the given ordering occurring in F . We say that a generalized AADD of form $c + bF$ satisfies the order if F satisfies the order.

Lemma 2.2. *Fix a variable ordering. For any non-constant function g mapping $\mathbb{B}^n \rightarrow \mathbb{R}$, there exists a unique generalized AADD G satisfying the given variable ordering such that for all $\rho \in \mathbb{B}^n$ we have $g(\rho) = Val(G, \rho)$.*

Both lemmas can be proved by straightforward induction on n . The second lemma shows that under a given variable ordering, generalized AADDs are canonical, i.e., two identical functions will always have identical AADD representations.

As an example of two AADDs with compact additive and multiplicative structure, Figure 2 shows portions of

¹Since normalized AADDs in grammar F are restricted to the range $[0, 1]$, we need the top-level positive affine transform of generalized AADDs in grammar G to allow for the representation of functions with arbitrary range.

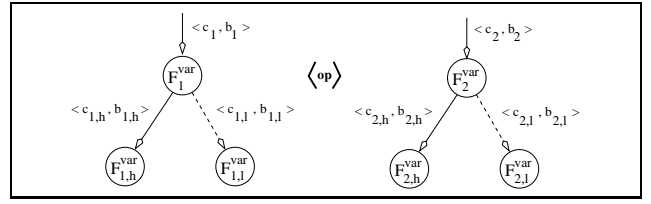


Figure 3: Two AADD nodes F_1 and F_2 and a binary operation op with the corresponding notation used in this paper.

the exponentially sized ADDs from Figure 1c represented by generalized AADDs of linear size.

We let op denote a binary operator on AADDs with possible operations being addition, subtraction, multiplication, division, min, and max denoted respectively as \oplus , \ominus , \otimes , \oslash , $\min(\cdot, \cdot)$, and $\max(\cdot, \cdot)$.

3 Related Work

There has been much related work in the formal verification literature that has attempted to tackle additive and multiplicative structure in representation of functions from $\mathbb{B}^n \rightarrow \mathbb{B}^m$. These include *BMDs, K*BMDs, EVBDDs, FEVBDDs, HDDs, and *PHDDs [3]. While space limitations prevent us from discussing all of the differences between AADDs and these data structures, we note two major differences: 1) These data structures have different canonical forms relying on GCD factorization and do not satisfy the invariant property of AADDs that internal nodes have range $[0, 1]$. 2) The fact that these data structures have integer terminals (\mathbb{B}^m) requires a rational or direct floating-point representation of values in \mathbb{R} . In our experience, this renders them unusable for probabilistic inference due to considerable numerical precision error and computational difficulties.

4 Algorithms

We now define AADD algorithms that are analogs of those given by Bryant [4]² except that they are extended to propagate the affine transform of the edge weights on recursion and to compute the normalization of the resulting node on return.

All algorithms rely on the helper function *getNode* given in Algorithm 1 that takes an unnormalized AADD node of the form $if(v, c_h + b_h F_h, c_l + b_l F_l)$ and returns the cached, generalized AADD node of the form $c_r + b_r F_r$.

4.1 Reduce

The *Reduce* algorithm given in Algorithm 2 takes an arbitrary ordered AADD, normalizes and caches the internal nodes, and returns the corresponding generalized AADD. One nice property of the *Reduce* algorithm is that one does not need to prespecify the structure that the AADD should exploit. If the represented function contains context-specific, additive, or multiplicative independence, the *Reduce* algorithm will compactly represent this structure uniquely and automatically.

²Bryant gives algorithms for binary decision diagrams (BDDs), but they are essentially identical to those for ADDs.

Algorithm 1: $GetNode(v, \langle c_l, b_l, F_l \rangle, \langle c_h, b_h, F_h \rangle) \rightarrow \langle c_r, b_r, F_r \rangle$

input : $v, \langle c_l, b_l, F_l \rangle, \langle c_h, b_h, F_h \rangle$: Var, offset, mult, and node id for low/high branches

output : $\langle c_r, b_r, F_r \rangle$: Return values for offset, multiplier, and canonical node id

begin

```

// If branches redundant, return child
if  $(c_l = c_h \wedge b_l = b_h \wedge F_l = F_h)$  then
   $\perp$  return  $\langle c_l, b_l, F_l \rangle$ ;
// Non-redundant so compute canonical form
 $r_{min} := \min(c_l, c_h)$ ;  $r_{max} := \max(c_l + b_l, c_h + b_h)$ ;
 $r_{range} := r_{max} - r_{min}$ ;
 $c_l := (c_l - r_{min}) / r_{range}$ ;  $c_h := (c_h - r_{min}) / r_{range}$ ;
 $b_l := b_l / r_{range}$ ;  $b_h := b_h / r_{range}$ ;
if  $(\langle v, \langle c_l, b_l, F_l \rangle, \langle c_h, b_h, F_h \rangle \rightarrow id$  is not in node cache) then
   $\perp$   $id :=$  currently unallocated id;
   $\perp$  insert  $\langle v, \langle c_l, b_l, F_l \rangle, \langle c_h, b_h, F_h \rangle \rightarrow id$  in cache;
// Return the cached, canonical node
return  $\langle r_{min}, r_{range}, id \rangle$ ;

```

end

4.2 Apply

The *Apply* routine given in Algorithm 3 takes two generalized AADD operands and an operation as given in Figure 3 and produces the resulting generalized AADD. Intuitively, the *Apply* algorithm produces the directed-acyclic graph (DAG) conforming to a joint-traversal of both operand DAGs and reduces it to canonical form on return. Modifications to the basic ADD *Apply* algorithm include the previously discussed propagation of the edge weight affine transformation on recursion and normalization on return, but also a number of pruning optimizations and a crucial scheme for canonically caching the results of *Apply* operations for optimal reuse.

Terminal computation and pruning

The function *ComputeResult* given in the top half of Table 1, determines if the result of a computation can be immediately computed without recursion. The first entry in this table is required for proper termination of the algorithm as it computes the result of an operation applied to terminal nodes. However, the other entries denote a number of pruning optimizations that immediately return a node without recursion. For example, given the operation $\langle 3 + 4F_1 \rangle \oplus \langle 5 + 6F_1 \rangle$, we can immediately return the result $\langle 8 + 10F_1 \rangle$.

Canonical caching

If the AADD *Apply* algorithm were to compute and cache the results of applying an operation directly to the operands, the algorithm would provably have the same time complexity as the ADD *Apply* algorithm. Yet, if we were to compute $\langle 0 + 1F_1 \rangle \oplus \langle 0 + 2F_2 \rangle$ and cache the result $\langle c_r + b_r F_r \rangle$, we could compute $\langle 5 + 2F_1 \rangle \oplus \langle 4 + 4F_2 \rangle = \langle (2c_r + 9) + 2b_r F_r \rangle$ without recursion.

Algorithm 2: $Reduce(\langle c, b, F \rangle) \rightarrow \langle c_r, b_r, F_r \rangle$

input : $\langle c, b, F \rangle$: Offset, multiplier, and node id

output : $\langle c_r, b_r, F_r \rangle$: Return values for offset, multiplier, and node id

begin

```

// Check for terminal node
if  $(F = 0)$  then
   $\perp$  return  $\langle c, 0, 0 \rangle$ ;
// Check reduce cache
if  $(F \rightarrow \langle c_r, b_r, F_r \rangle$  is not in reduce cache) then
  // Not in cache, so recurse
   $\langle c_h, b_h, F_h \rangle := Reduce(c_h, b_h, F_h)$ ;
   $\langle c_l, b_l, F_l \rangle := Reduce(c_l, b_l, F_l)$ ;
  // Retrieve canonical form
   $\langle c_r, b_r, F_r \rangle := GetNode(F^{var}, \langle c_l, b_l, F_l \rangle, \langle c_h, b_h, F_h \rangle)$ ;
  // Put in cache
   $\perp$  insert  $F \rightarrow \langle c_r, b_r, F_r \rangle$  in reduce cache;
// Return canonical reduced node
return  $\langle b \cdot c_r + c, b \cdot b_r, F_r \rangle$ ;

```

end

This suggests a canonical caching scheme that normalizes all cache entries to increase the probability of a cache hit. The actual result can then be easily computed from the cached result by reversing the normalization. This ensures optimal reuse of the *Apply* operations cache and can lead to an exponential reduction in running time over the non-canonical caching version.

We introduce two additional functions to perform this caching: *GetNormCacheKey* to compute the canonical cache key, and *ModifyResult* to reverse the normalization in order to compute the actual result. These algorithms are summarized in the lower half of Table 1.

4.3 Other Operations

While space limitations prevent us from covering all of the operations that can be performed (efficiently) on AADDs, we briefly summarize some of them here:

- **min and max computation:** The min and max of a generalized AADD node $\langle c + bF \rangle$ are respectively c and $c + b$ due to $[0, 1]$ normalization of F .
- **Restriction:** The restriction of a variable x_i in a function to either *true* or *false* (i.e. $F|_{x_i=T/F}$) can be computed by returning the proper branch of nodes containing that test variable during the *Reduce* operation.
- **Sum out/marginalization:** A variable x_i can be summed (or marginalized) out of a function F simply by computing the sum of the restricted functions (i.e. $F|_{x_i=T} \oplus F|_{x_i=F}$).
- **Variable reordering:** A generalization of Rudell's ADD variable reordering algorithm [5] that recomputes edge weights of reordered nodes can be applied to AADDs without loss of efficiency.

5 Theoretical results

Here we present two fundamental results for AADDs:

Theorem 5.1. *The time and space performance of Reduce and Apply for AADDs is within a multiplicative constant of that of ADDs in the worst case.*

Proof Sketch. Under the same variable ordering, an ADD is equivalent to a non-canonical AADD with fixed edge weights $c = 0, b = 1$. Thus the ADD *Reduce* and *Apply* algorithms can be seen as analogs of the AADD algorithms without the additional constant-time and constant-space overhead of normalization. Since normalization can only increase the number of *Reduce* and *Apply* cache hits and reduce the number of cached nodes, it is clear that an AADD must generate equal or fewer *Reduce* and *Apply* calls and have equal or fewer cached nodes than the corresponding ADD. This allows us to conclude that in the worst case where the AADD generates as many *Reduce* and *Apply* calls and cache hits as the ADD, the AADD is still within a multiplicative constant of the time and space required by the ADD.

Theorem 5.2. *There exist functions F_1 and F_2 and an operator op such that the running time and space performance of $Apply(F_1, F_2, op)$ for AADDs can be linear in the number of variables when the corresponding ADD operations are exponential in the number of variables.*

Proof Sketch. Two functions and *Apply* operation examples where this holds true are $\sum_{i=1}^n 2^i x_i \oplus \sum_{i=1}^n 2^i x_i$ and $\prod_{i=1}^n \gamma^{2^i x_i} \otimes \prod_{i=1}^n \gamma^{2^i x_i}$. (Simple examples of these operands for ADDs and AADDs are given respectively in Figures 1c and 2.) Because the results of these operations have a number of terminal values exponential in n , the ADD operations must require time and space exponential in n . On the other hand, it is known that the operands can be represented in linear-sized AADDs, and once the AADD *Apply* operation computes \oplus or \otimes for all n low branches of operand nodes with matching variables, computation of the corresponding high branches will yield cache hits due to the normalized caching scheme given in Table 1. This results in n cached nodes and $2n$ *Apply* calls for the AADD operations on these functions.

6 Empirical results

6.1 Basic operations

Figure 4 demonstrates the relative time and space performance of ADDs, AADDs, and tables on a number of basic operations. These verify the exponential to linear space and time reductions proved in Theorem 5.2.

6.2 Bayes nets

For Bayes nets (BNs), we simply evaluate the variable elimination algorithm [6] (under a fixed, greedy tree-width minimizing variable ordering) with the conditional probability tables (CPT) P_i and corresponding operations replaced with those for tables, ADDs, and AADDs:

$$\sum_{x_1 \dots x_i} \prod_{P_1 \dots P_j} P_1(x_k | x_l, \dots, x_m) \dots P_j(x_n | x_p, \dots, x_q)$$

Algorithm 3: $Apply(\langle c_1, b_1, F_1 \rangle, \langle c_2, b_2, F_2 \rangle, op) \rightarrow \langle c_r, b_r, F_r \rangle$

```

input :  $\langle c_1, b_1, F_1 \rangle, \langle c_2, b_2, F_2 \rangle, op$  : Nodes and op
output :  $\langle c_r, b_r, F_r \rangle$  : Generalized node to return
begin
  // Check if result can be immediately computed
  if ( $ComputeResult(\langle c_1, b_1, F_1 \rangle, \langle c_2, b_2, F_2 \rangle, op) \rightarrow \langle c_r, b_r, F_r \rangle$  is not null) then
     $\perp$  return  $\langle c_r, b_r, F_r \rangle$ ;
  // Get normalized key and check apply cache
   $\langle \langle c'_1, b'_1 \rangle, \langle c'_2, b'_2 \rangle \rangle :=$ 
     $GetNormCacheKey(\langle c_1, b_1, F_1 \rangle, \langle c_2, b_2, F_2 \rangle, op)$ ;
  if ( $\langle \langle c'_1, b'_1 \rangle, \langle c'_2, b'_2 \rangle \rangle \rightarrow \langle c_r, b_r, F_r \rangle$  is not in apply cache) then
    // Not terminal, so recurse
    if ( $F_1$  is a non-terminal node) then
      if ( $F_2$  is a non-terminal node) then
        if ( $F_1^{var}$  comes before  $F_2^{var}$ ) then
          var :=  $F_1^{var}$ ;
        else
           $\perp$  var :=  $F_2^{var}$ ;
        else
           $\perp$  var :=  $F_1^{var}$ ;
      else
         $\perp$  var :=  $F_2^{var}$ ;
    // Propagate affine transform to branches
    if ( $F_1$  is non-terminal  $\wedge$  var =  $F_1^{var}$ ) then
       $F_l^{v1} := F_{1,l}; F_h^{v1} := F_{1,h};$ 
       $c_l^{v1} := c'_1 + b'_1 \cdot c_{1,l}; c_h^{v1} := c'_1 + b'_1 \cdot c_{1,h};$ 
       $b_l^{v1} := b'_1 \cdot b_{1,l}; b_h^{v1} := b'_1 \cdot b_{1,h};$ 
    else
       $\perp F_{l/h}^{v1} := F_1; c_{l/h}^{v1} := c'_1; b_{l/h}^{v1} := b'_1;$ 
    if ( $F_2$  is non-terminal  $\wedge$  var =  $F_2^{var}$ ) then
       $F_l^{v2} := F_{2,l}; F_h^{v2} := F_{2,h};$ 
       $c_l^{v2} := c'_2 + b'_2 \cdot c_{2,l}; c_h^{v2} := c'_2 + b'_2 \cdot c_{2,h};$ 
       $b_l^{v2} := b'_2 \cdot b_{2,l}; b_h^{v2} := b'_2 \cdot b_{2,h};$ 
    else
       $\perp F_{l/h}^{v2} := F_2; c_{l/h}^{v2} := c'_2; b_{l/h}^{v2} := b'_2;$ 
    // Recurse and get cached result
     $\langle c_l, b_l, F_l \rangle := Apply(\langle c_l^{v1}, b_l^{v1}, F_l^{v1} \rangle, \langle c_l^{v2}, b_l^{v2}, F_l^{v2} \rangle, op)$ ;
     $\langle c_h, b_h, F_h \rangle := Apply(\langle c_h^{v1}, b_h^{v1}, F_h^{v1} \rangle, \langle c_h^{v2}, b_h^{v2}, F_h^{v2} \rangle, op)$ ;
     $\langle c_r, b_r, F_r \rangle := GetNode(var, \langle c_l, b_l, F_l \rangle, \langle c_h, b_h, F_h \rangle)$ ;
    // Put result in apply cache and return
    insert  $\langle c'_1, b'_1, F_1, c'_2, b'_2, F_2, op \rangle \rightarrow \langle c_r, b_r, F_r \rangle$ 
    into apply cache;
  return  $ModifyResult(\langle c_r, b_r, F_r \rangle)$ ;
end

```

Table 2 shows the total number of table entries/nodes required to represent the original network and the total running time of 100 random queries (each consisting of one query variable and one evidence variable) for a number of publicly available BNs (<http://www.cs.huji.ac.il/labs/compbio/Repository>) and two *noisy-or* and *noisy-max* CPTs with 15 parent nodes.

$ComputeResult(\langle c_1, b_1, F_1 \rangle, \langle c_2, b_2, F_2 \rangle, op) \rightarrow \langle c_r, b_r, F_r \rangle$	
Operation and Conditions	Return Value
$\langle c_1 + b_1 F_1 \rangle \langle op \rangle \langle c_2 + b_2 F_2 \rangle; F_1 = F_2 = 0$	$\langle (c_1 \langle op \rangle c_2) + 0 \cdot 0 \rangle$
$\max(\langle c_1 + b_1 F_1 \rangle, \langle c_2 + b_2 F_2 \rangle); c_1 + b_1 \leq c_2$	$\langle c_2 + b_2 F_2 \rangle$
$\max(\langle c_1 + b_1 F_1 \rangle, \langle c_2 + b_2 F_2 \rangle); c_2 + b_2 \leq c_1$	$\langle c_1 + b_1 F_1 \rangle$
$\langle c_1 + b_1 F_1 \rangle \oplus \langle c_2 + b_2 F_2 \rangle; F_1 = F_2$	$\langle (c_1 + c_2) + (b_1 + b_2) F_1 \rangle$
$\max(\langle c_1 + b_1 F_1 \rangle, \langle c_2 + b_2 F_1 \rangle); F_1 = F_2,$ $(c_1 > c_2 \wedge b_1 > b_2) \vee (c_2 > c_1 \wedge b_2 > b_1)$	$c_1 \geq c_2 \wedge b_1 \geq b_2 : \langle c_1 + b_1 F_1 \rangle$ $c_2 \geq c_1 \wedge b_2 \geq b_1 : \langle c_2 + b_2 F_1 \rangle$
Note: for all max operations above, return opposite for min	
$\langle c_1 + b_1 F_1 \rangle \langle op \rangle \langle c_2 + b_2 F_2 \rangle; F_2 = 0, op \in \{\oplus, \ominus\}$	$\langle (c_1 \langle op \rangle c_2) + b_1 F_1 \rangle$
$\langle c_1 + b_1 F_1 \rangle \langle op \rangle \langle c_2 + b_2 F_2 \rangle; F_2 = 0, c_2 \geq 0, op \in \{\otimes, \oslash\}$	$\langle (c_1 \langle op \rangle c_2) + (b_1 \langle op \rangle c_2) F_1 \rangle$
Note: above two operations can be modified for $F_1 = 0$ when $op \in \{\oplus, \otimes\}$	
other	<i>null</i>

$GetNormCacheKey(\langle c_1, b_1, F_1 \rangle, \langle c_2, b_2, F_2 \rangle, op) \rightarrow \langle \langle c'_1, b'_1 \rangle \langle c'_2, b'_2 \rangle \rangle$ and $ModifyResult(\langle c_r, b_r, F_r \rangle) \rightarrow \langle c'_r, b'_r, F'_r \rangle$		
Operation and Conditions	Normalized Cache Key and Computation	Result Modification
$\langle c_1 + b_1 F_1 \rangle \oplus \langle c_2 + b_2 F_2 \rangle; F_1 \neq 0$	$\langle c_r + b_r F_r \rangle = \langle 0 + 1 F_1 \rangle \oplus \langle 0 + (b_2/b_1) F_2 \rangle$	$\langle (c_1 + c_2 + b_1 c_r) + b_1 b_r F_r \rangle$
$\langle c_1 + b_1 F_1 \rangle \ominus \langle c_2 + b_2 F_2 \rangle; F_1 \neq 0$	$\langle c_r + b_r F_r \rangle = \langle 0 + 1 F_1 \rangle \ominus \langle 0 + (b_2/b_1) F_2 \rangle$	$\langle (c_1 - c_2 + b_1 c_r) + b_1 b_r F_r \rangle$
$\langle c_1 + b_1 F_1 \rangle \otimes \langle c_2 + b_2 F_2 \rangle; F_1 \neq 0$	$\langle c_r + b_r F_r \rangle = \langle (c_1/b_1) + F_1 \rangle \otimes \langle (c_2/b_2) + F_2 \rangle$	$\langle b_1 b_2 c_r + b_1 b_2 b_r F_r \rangle$
$\langle c_1 + b_1 F_1 \rangle \oslash \langle c_2 + b_2 F_2 \rangle; F_1 \neq 0$	$\langle c_r + b_r F_r \rangle = \langle (c_1/b_1) + F_1 \rangle \oslash \langle (c_2/b_2) + F_2 \rangle$	$\langle (b_1/b_2) c_r + (b_1/b_2) b_r F_r \rangle$
$\max(\langle c_1 + b_1 F_1 \rangle, \langle c_2 + b_2 F_2 \rangle);$ $F_1 \neq 0$, Note: same for min	$\langle c_r + b_r F_r \rangle = \max(\langle 0 + 1 F_1 \rangle, \langle (c_2 - c_1)/b_1 + (b_2/b_1) F_2 \rangle)$	$\langle (c_1 + b_1 c_r) + b_1 b_r F_r \rangle$
any $\langle op \rangle$ not matching above: $\langle c_1 + b_1 F_1 \rangle \langle op \rangle \langle c_2 + b_2 F_2 \rangle$	$\langle c_r + b_r F_r \rangle = \langle c_1 + b_1 F_1 \rangle \langle op \rangle \langle c_2 + b_2 F_2 \rangle$	$\langle c_r + b_r F_r \rangle$

Table 1: Input and output summaries of the *ComputeResult*, *GetNormCacheKey*, and *ModifyResult* routines.

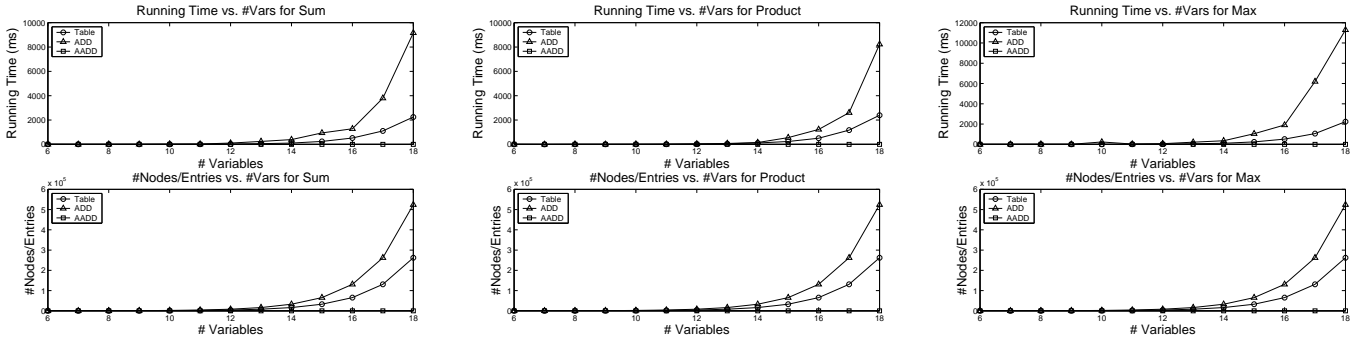


Figure 4: A comparison of running time (top) and table entries/nodes (bottom) for Tables, ADDs and AADDs under various *Apply* operations. Left to right: $(\sum_i 2^i x_i) \oplus (\sum_i 2^i x_i)$, $(\gamma^{\sum_i 2^i x_i}) \otimes (\gamma^{\sum_i 2^i x_i})$, $\max(\sum_i 2^i x_i, \sum_i 2^i x_i)$.

Note that the intermediate probability tables were often too large for the tables or ADDs, but not the AADDs, indicating that the AADD was able to exploit additive or multiplicative structure in these cases. Also, the AADD yields an exponential to linear reduction on the *Noisy-Or-15* problem and a considerable computational speedup on the *Noisy-Max-15* problem by exploiting the additive and multiplicative structure inherent in these special CPTs [7].

6.3 Markov decision processes

For MDPs, we simply evaluate the value iteration algorithm [8] for a tabular representation and its extension for decision diagrams [9] to factored MDPs from the

SysAdmin domain [10].³ Here we simply substitute the tabular representation, ADD, and AADD for the reward function R , value function V , and transition function T in the following value iteration update:

$$V^{n+1}(s_1, \dots, s_i) = R(s_1, \dots, s_i) + \gamma \max_a \sum_{s'_1, \dots, s'_i} T(s'_1, \dots, s'_i | s_1, \dots, s_i, a) V^n(s'_1, \dots, s'_i)$$

³In these domains, there is a network of computers whose operational status (i.e., running or crashed) at each time step is a stochastic function of the previous status of computers with incoming connections. At each time step the task of the system administrator is to choose the computer to reboot in order to maximize the up-time of all computers over a discounted reward horizon.

Bayes Net	Table		ADD		AADD	
	# Table Entries	Running Time	# ADD Nodes	Running Time	# AADD Nodes	Running Time
Alarm	1,192	2.97 s	689	2.42 s	405	1.26 s
Barley	470,294	<i>EML*</i>	139,856	<i>EML*</i>	60,809	207 m
Carmo	636	0.58	955	0.57 s	360	0.49 s
Hailfinder	9045	26.4 s	4511	9.6 s	2538	2.7 s
Insurance	2104	278 s	1596	116 s	775	37 s
Noisy-Or-15	65566	27.5 s	125356	50.2 s	1066	0.7 s
Noisy-Max-15	131102	33.4 s	202148	42.5 s	40994	5.8 s

Table 2: Number of table entries/nodes in the original network and variable elimination running times using tabular, ADD, and AADD representations for inference in various Bayes nets. * *EML* denotes that a query exceeded the 1Gb memory limit.

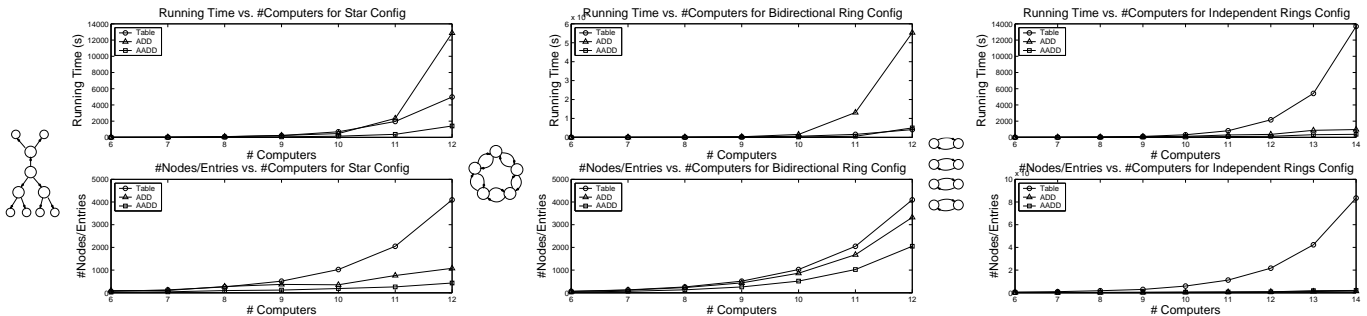


Figure 5: MDP value iteration running times (top) and number of entries/nodes (bottom) in the final value function using tabular, ADD, and AADD representations for various network configurations in the *SysAdmin* problem.

Figure 5 shows the relative performance of value iteration until convergence within 0.01 of the optimal value for networks in a star, bidirectional, and independent ring configuration. While the reward and transition dynamics in the *SysAdmin* problem have considerable additive structure, we note that the exponential size of the AADD (as for all representations) indicates that little additive structure survives in the *exact* value function. Nonetheless, the AADD algorithm still manages to take considerable advantage of the additive structure during computations and thus performs comparably or exponentially better than ADDs and tables.

7 Concluding Remarks

We have presented the AADD and have proved that its worst-case time and space performance are within a multiplicative constant of that of ADDs, but can be linear in the number of variables in cases where ADDs are exponential in the number of variables. And we have provided an empirical comparison of tabular, ADD, and AADD representations used in Bayes net and MDP inference algorithms, concluding that AADDs perform at least as well as the other two representations, and often yield an exponential time and space improvement over both.

In conclusion, we should emphasize that we have not set out to show that variable elimination and value iteration with AADDs are the best Bayes net and MDP inference algorithms available – many other approaches propose to handle similar structure efficiently via special-purpose problem-structure or algorithm modifications. Rather, we intended to show that transparently substituting AADDs in two diverse probabilistic inference al-

gorithms that used tables or ADDs could yield exponential performance improvements over both by exploiting context-specific, additive, and multiplicative structure in the underlying representation. These results suggest that the AADD is likely to yield exponential time and space improvements for a variety of probabilistic inference algorithms that currently use tables or ADDs, or at the very least, a gross simplification of these algorithms.

References

- [1] R.I. Bahar, E.A. Frohm, C.M. Gaona, G.D. Hachtel, E. Macii, A. Pardo, F. Somenzi: Algebraic Decision Diagrams and Their Applications. ICCAD. (1993)
- [2] Boutilier, C., Friedman, N., Goldszmidt, M., Koller, D.: Context-specific independence in Bayesian networks. UAI. (1996)
- [3] Drechsler, R., Sieling, D.: BDDs in theory and practice. Software Tools for Technology Transfer. (2001)
- [4] Bryant, R.E.: Graph-based algorithms for Boolean function manipulation. IEEE Trans. on Computers. (1986)
- [5] Rudell, R.: Dynamic variable ordering for ordered binary decision diagrams. ICCAD. (1993)
- [6] Zhang, N.L., Poole, D.: A simple approach to bayesian network computations. CCAI. (1994)
- [7] Takikawa, M., D'Ambrosio, B.: Multiplicative factorization of noisy-max. UAI. (1999)
- [8] Bellman, R.E.: Dynamic Programming. Princeton University Press, Princeton, NJ (1957)
- [9] Hoey, J., St-Aubin, R., Hu, A., Boutilier, C.: SPUDD: Stochastic planning using decision diagrams. UAI. (1999)
- [10] Guestrin, C., Koller, D., Parr, R.: Max-norm projections for factored MDPs. IJCAI. (2001)