

## بررسی تطبیقی الگوریتم‌های ساخت هیستوگرام‌های درخت R و مقایسه آن با

### مجموعه درخت‌های $R^+$ و $R^*$

محبوبه شمس، مرضیه دوستی

عضو هیئت علمی دانشکده برق و کامپیوتر دانشگاه صنعتی قم

## A comparative study of algorithms to build the tree histograms comparison with the R and R + and R \* of trees

Mahboobeh shamsi, marzie dousti

Faculty members of electrical and computer engineering university of qom

Email: shamsi@qut.ac.ir

Email: Dousti@gmail.com

### چکیده

امروزه طبقه‌بندی اطلاعات برای استفاده مناسب از آنها به هنگام نیاز امری بسیار مهم و حیاتی است، در این راستا روش‌های مختلفی برای طبقه‌بندی معرفی شده است. که هر یک به نوبه خود بهبودی را در راستای بهینه‌سازی اکتشاف و طبقه‌بندی داده بوجود آورده‌اند، یکی از روش‌های مهم در این مقوله، استفاده از شاخص‌های فضایی است. نظریه اصلی شاخص فضایی که مهمترین نظریه اساسی پرس و جو می‌باشد، محاسبه نظریه‌ی نزدیکی است. معروف‌ترین ساختار شاخص، درخت R می‌باشد. استراتژی اصلی درخت R جمع آوری نودهای فضایی چند بعدی با مستطیل حداقل محدوده (MBR) می‌باشد که کوچکترین نود فضایی داخلی مستطیل باشد. بعد از شاخص بندی، بهینه سازی بازیابی اطلاعات در پایگاه داده‌ی فضایی مهم است [۲]. از اینرو جهت بررسی بیشتر موضوع، در این مقاله، شاخص فضایی جدید که متعلق به خانواده  $R^+$  tree می‌باشد و مفهوم همپوشانی گره‌ها را حفظ می‌کند و  $R^{++}$  tree نام دارد، بررسی شده است. نتایج این مقاله نشان می‌دهد که  $R^{++}$  tree در دامنه‌ی پرس و جو، پرس و جو KNN و پرس و جو Top-k بسیار کارآمدتر از  $R^*$  tree می‌باشد.

کلمات کلیدی:

شاخص فضایی، هیستوگرام‌های ساخت درخت‌های R، درخت‌های R

### ۱. مقدمه

هیستوگرام‌ها، ساختارهای مهمی هستند که در ابتدا در سیستم‌های پایگاه داده به منظور تخمین انتخابی پرس و جوها مورد استفاده قرار گرفتند. آنها برای کسب پاسخ تقریبی سریع برای پرس و جوهای مترکم مهم استفاده می‌شوند. بیشتر آنها متدهای مبتنی بر شبکه‌ای هستند که تنها برای دیتا پوینت دو بعدی قابل اجرا هستند. مشکل طراحی بهینه هیستوگرام‌های چندبعدی به NP-hard مشهور هستند [۳].

وضعیت فعلی بیانگر این است که حین اجرای یک پرس و جو در نظام مدیریت یک پایگاه داده‌ی فضایی (SDBMS)، بهینه‌سازی پرس و جو تمام طرح‌های موجود در ارزیابی پرس و جو را به وجود می‌آورد. تمام این طرح‌ها در بازده نهایی مشابه اند اما در هزینه اجرا و زمان اجرا با یکدیگر فرق می‌کنند [۲].

فلسفه تخصیص باکت‌های هیستوگرام اختصاص آنها به زیرفضاهایی که بدرستی خوشه‌های اشیا را تشخیص داده اند می‌باشد. بنابراین، ابتدا یک روش برای پیدا کردن مرکز خوشه‌های اشیای پیشنهاد می‌شود. سپس، یک الگوریتم برای ساخت باکت‌های هیستوگرام از این مراکز پیشنهاد می‌گردد. این باکت‌ها از مراکز خوشه‌ها، مقداره‌ی اولیه می‌شوند، سپس برای پوشش خوشه‌ها گسترش می‌یابند. بهترین طرح توسعه براساس مفهوم افزایش چولگی انتخاب شده است [۶].

## ۲. تاریخچه

یکی از اولین متدها که برای دیتای چندبعدی پیشنهاد شد htree می‌باشد. این متد پارتیشن بندی ناهمپوشی فضای چند بعدی را بر اساس فرکانس به عنوان پارامتر منشا ایجاد می‌کند [۲] نویسندگان دو استراتژی را پیشنهاد کرده اند: (۱) متغیر بنیادی بصورت زیر عمل می‌کند: در فاز اول، الگوریتم شبکه معینی را محاسبه می‌کند و تعداد اشیای فضایی زیر گروه را برای هرخانه تعریف می‌کند. براساس شبکه محاسبه شده تقسیم بندی ساخت دوتایی بازگشتی (BSP) برای محاسبه هیستوگرام استفاده می‌شود.

(۲) برای کاهش تاثیر استراتژی دوم ساختاری یعنی Minskev-Progressive-Refinement شبکه‌هایی با تفکیک‌های مختلف مورد استفاده قرار می‌گیرند. هر تفکیک شبکه‌ای برای ساخت بخش‌های یکسان منحنی‌های پیوند هیستوگرام استفاده می‌شود [۲].

MinSkew یک روش شناخته شده ساخت هیستوگرام برای داده‌های فضایی است. در ابتدا، MinSkew مجموعه داده‌های اصلی را با استفاده از یک شبکه یکنواخت تخمین می‌زند. سپس، با یک باکت تک متشکل از تمام اشیاء داده شروع می‌شود، برای هر باکت، انحراف مکانی آن و نقطه تقسیم در طول ابعادش که حداکثر کاهش در انحراف مکانی را تولید خواهد کرد، محاسبه می‌کند و بعد، MinSkew باکت‌هایی که تقسیم آنها منجر به بزرگترین کاهش در چوله مکانی می‌شود، را بر می‌دارد. این باکت را به دو باکت فرزند تقسیم می‌کند، و داده را از باکت قدیمی به باکت جدید اختصاص می‌دهد. پس از اتمام MinSkew، هیستوگرام ساخته شده مجموعه‌ای از باکت غیر متداخل است [۶].

سادگی می‌توان از الگوریتم ساخت باکت Bichromatic برای دریافت یک نسخه جدید از هر باکت ساخته شده استفاده کرد. مجموعه باکت‌های Bichromatic سپس به عنوان هیستوگرام نهایی گزارش می‌شود. این استراتژی، به عنوان مثال، تبدیل هر باکت به Bichromatic پس از اتمام فرآیند ساخت هیستوگرام، می‌تواند برای هر روشی که در آن باکت ساخته شده با یکدیگر همپوشانی ندارند، بکار گرفته شود [۶].

STHist یک روش ساخت هیستوگرام برای داده‌های جغرافیایی گرافیکی دو یا سه بعدی است. با توجه به مجموعه داده همراه با فضای داده‌ها، STHist ابتدا تمام فضای داده‌ها را به تعدادی از بخش‌های داده پارتیشن بندی می‌کند. سپس، برای هر بخش، STHist به صورت بازگشتی نقاط متراکم که به باکت هیستوگرام تبدیل شده است را تشخیص می‌دهد. در اینجا، نقاط متراکم یک منطقه داده است که شرایط خاصی در فرکانس جسم، شکل و اندازه را برآورده می‌کند. همه باکت‌ها در یک بخش داده که یک درخت باکت سازمان یافته است شناسایی شدند. با توجه به روش ساخت باکت Bichromatic، ما می‌توانیم از الگوریتم ساخت باکت Bichromatic برای تبدیل هر باکت ساخته شده توسط STHist در طول فرآیند ساخت هیستوگرام به نسخه Bichromatic استفاده کنیم. به عبارت دیگر، برای هر بخش داده‌ها، پس از اینکه باکت ریشه ایجاد می‌شود، الگوریتم ساخت باکت Bichromatic برای این باکت ریشه استفاده می‌شود. سپس، STHist نقاط داخل باکت ریشه بهبود یافته را تشخیص می‌دهد و این نقاط را به باکت فرزند تبدیل می‌کند. برای هر یک از این باکت‌های فرزند، ما الگوریتم ساخت باکت Bichromatic را برای دریافت نسخه بهبود یافته اعمال می‌کنیم. این روند تا زمانی که هیچ باکت جدیدی ساخته نشده باشد ادامه می‌یابد [۶].

STHist-c روش بهبود یافته STHist می‌باشد که همان چهارچوب STHist را بکار می‌برد. روش STHist شرایط سفت و سخت را بر روی شکل و اندازه برای تشخیص مناطق تراکم بکار می‌برد. این دلیل عمدتاً برای یک راه مرتبط آسان برای یافتن مناطق تراکم می‌باشد. بهر حال بخاطر این شرایط سفت و سخت، روش STHist گاهی اوقات موفق به تشخیص

خوشه‌های اشیاء در مجموعه داده ها نمی‌شود. برای ساخت نمودار هیستوگرام دقیق، STHist-c باکت‌ها را در مکان خوشه‌های داده اختصاص می‌دهد. ابتدا یک الگوریتم خوشه بندی و دو روش آماری، برای پیدا کردن تعداد و مکان خوشه‌های اشیاء با هم ترکیب می‌شوند. سپس یک الگوریتم جدید برای ساخت باکت‌ها از مرکز این خوشه‌ها ارائه می‌شود. مناطق باکت‌ها به تدریج از مرکز خوشه به تمام جهات گسترش می‌یابند. در زمان گسترش یک باکت، از میان بسیاری از گزینه‌های احتمالی گسترش، بهترین روش بر اساس مفهوم جدیدی از افزایش چولگی انتخاب می‌شود [۶].

### ۳. آزمایشات

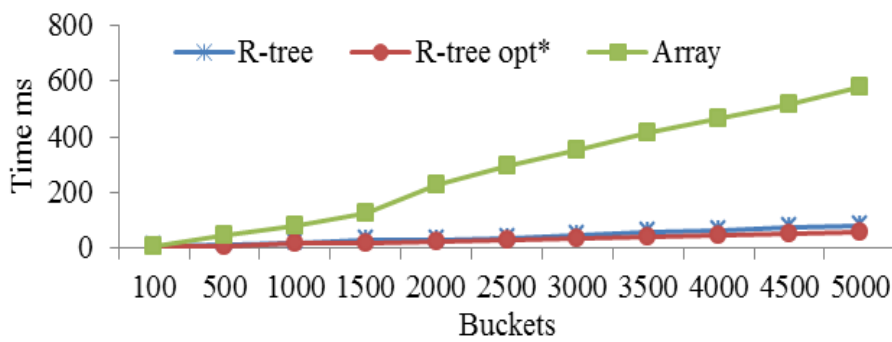
در این مقاله، اجرای هیستوگرام‌های مختلف را مورد بررسی و مقایسه قرار گرفته شده است. از MinSkew بعنوان متد مرجع استفاده شده که بسیار خوب عمل می‌کند، دومی را با عنوان MinSkewProg معرفی شده است. برای  $d=2$  شبکه ای با خانه‌های  $2^4$  مورد استفاده قرار داده شد. برای  $d=3$  خانه‌های  $2^{15}$  برای MinSkew و چهار شبکه با خانه‌های  $2^6$ ،  $2^9$ ،  $2^{12}$ ،  $2^{15}$  برای MinSkewProg مورد استفاده قرار داده شده است. دیگر متد‌ها که مورد آزمایش قرار گرفتند در جدول شماره ۱ لیست شده اند [۲].

جدول ۱: (بررسی تابع هزینه متدهای مطالعه شده) [۲]

توضیحات	تابع هزینه
حجم MBR	$C_v$
گسترش $C_v$ با میانگین، طولهای طرف پرس و جو	$C_{QP}$
معیار k-Uniformity	$C_{RK}$
Spatial skew of MBR	$C_{SK}$
توضیحات	هیستوگرام
MinSkew، شبکه ثابت	MinSkew
MinSkew,Prog، پالایش	MinSkewProg
rkHist with $\alpha=0.1$	rkHist
پارتیشن بندی سایز ثابت، منحنی Hilbert	Rtree
$C_v$ منحنی Hilbert	R-V
$C_{QP}$ منحنی Hilbert	R-VQP
$C_{RK}$ منحنی Hilbert	R-RK
$C_{SK}$ منحنی Hilbert	R-SK
جنگل STHist	FST

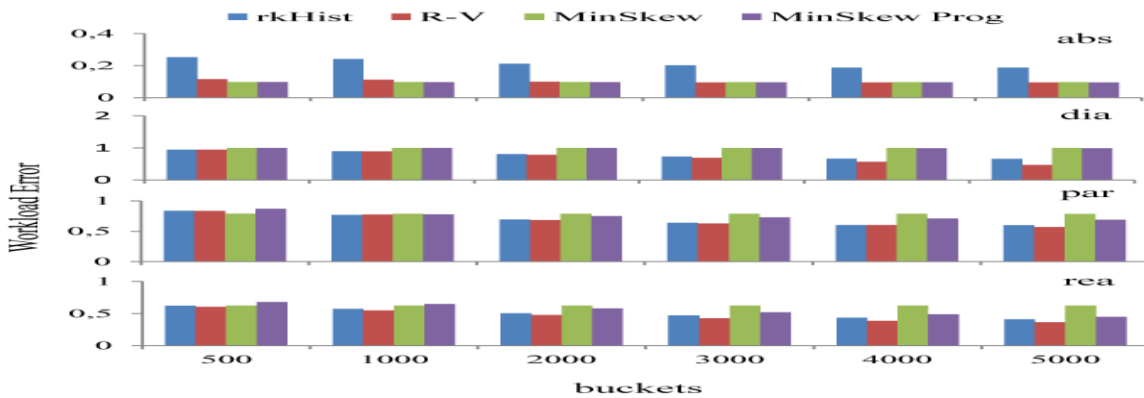
### ۴. زمان ساخت و تخمین

برای کاهش زمان تقریبی هیستوگرام می‌تواند بعنوان حافظه‌ی اصلی R-tree ارائه شود. شکل ۴ عملکرد اندازه‌ی منحنی‌های پیوندی و زمان مجموع حجم کار برای ارائه هیستوگرام را نشان می‌دهد، دوتای اول R-tree هستند. سومی آرایه‌ای برای منحنی‌های پیوند است. برای R-tree تنظیمات حافظه‌ی اصلی مورد استفاده قرار داده شده و گنجایش خروجی به ۱۲ مدخل برای هر گروه تنظیم شده است (دوباره ذکر شده بهترین مجموعه در آزمایشات بود). علاوه بر آن یک R-tree با استفاده از منحنی‌های پیوند هیستوگرام با یک متد پارتیشن بندی otp و  $C_v$  بعنوان تابع هزینه‌ای درست شده است [۲].



شکل ۱، [۲]

شکل ۲ نتایج متدهای R-tree را در مقایسه با استراتژی پارتیشن بندی اندازه ثابت برای دیتاپونیت  $d=2$  را نشان می‌دهد. مشاهده شد که متدهایی که از چهارچوب پارتیشن بندی بهینه بر اساس نوع استفاده می‌کنند دقت بهتری را نسبت R-tree از خود نشان می‌دهند.



شکل ۲، [۲]

یک راه حل ممکن پارتیشن بندی توزیع دیتایی بر اساس اندازه و شکل هدف و ساخت هیستوگرام یا شاخص گذاری بطور مستقل برای هر طبقه بندی می‌باشد [۲].

## ۵. تفاوت‌های درخت‌های $R^+$ و $R^*$

درخت  $R^+$  یک روش برای دنبال کردن داده با استفاده از یک محل اغلب مختصات  $(x,y)$ ، و اغلب برای مکان‌های بر روی سطح زمین است. جستجو در یک عدد یک مشکل حل شده است. جستجو در دو یا بیشتر، و درخواست برای مکان‌هایی که در نزدیکی هر دو مختصات  $x$  و  $y$ ، نیاز به الگوریتم‌های craftier دارد. اساساً، درخت  $R^+$  ساختمان داده‌ی یک نوع درخت  $R$ ، مورد استفاده برای اطلاعات مکانی شاخص گذاری است [۵].

### ۱.۵ تفاوت Rtree و $R^+$ tree:

$R^+$ tree یک مصالحه بین Rtree و kd-tree: آن‌ها از تداخل گره‌های داخلی با وارد کردن یک شیء به برگ‌های متعدد در صورت لزوم اجتناب می‌کنند. همپوشانی کل منطقه است که در دو یا چند گره موجود است. همپوشانی حداقل مجموعه‌ای از مسیرهای جستجو به برگ (حتی بیشتر برای زمان دسترسی از حداقل پوشش بحرانی) را کاهش می‌دهد. جستجوی کارآمد نیاز به حداقل پوشش و همپوشانی دارد [۵].



درخت  $R^+$  و درخت  $R$  در موارد زیر با هم تفاوت دارند:

گره‌ها حداقل نیمی پر تضمین نمی‌شود. ورودی‌های هر یک از گره‌های داخلی همپوشانی ندارند. یک شناسه شیئی می‌تواند در بیش از یک گره برگ ذخیره شود [۵].  
مزایا:

از آنجا که گره‌ها با یکدیگر همپوشانی ندارند. مزایای کارایی پرس و جوی نقطه از مناطق فضایی توسط حداکثر یک گره پوشیده شده است. یک مسیر واحد دنبال می‌شود و گره‌های کمتری نسبت به  $R_{tree}$  مشاهده می‌شوند [۵].  
معایب:

از آنجا که مستطیل‌ها تکرار می‌شوند، یک درخت  $R^+_{tree}$  می‌تواند بزرگتر از درخت  $R_{tree}$  بر روی مجموعه داده یکسان ساخته شود. ساخت و نگهداری  $R^+_{tree}$  پیچیده‌تر از ساخت و ساز و تعمیر و نگهداری دیگر انواع درخت  $R$  است [۵].  
درخت  $R^*$  یک نوع از درختان  $R$  هستند که برای شاخص گذاری اطلاعات مکانی استفاده می‌شوند. هزینه ساخت و ساز  $R^*_{tree}$  کمی بالاتر از استاندارد درختان  $R$  است. از آنجا که داده‌ها ممکن است نیاز به جاگذاری دوباره داشته باشند، اما درخت نتیجه معمولاً عملکرد پرس و جو بهتر خواهد داشت. مانند درخت استاندارد  $R$ ، می‌تواند در هر دو نقطه و داده‌های مکانی ذخیره شود [۵].

### ۲.۵. تفاوت بین $R_{tree}$ و $R^*_{tree}$ :

به حداقل رساندن همپوشانی و پوشش هر دو به کارایی درخت  $R$  بستگی دارد. همپوشانی بدان معنی است که، در پرس و جوی داده‌ها و یا درج، بیش از یک شاخه از درخت نیاز به گسترش دارد. پوشش حداقلی باعث بهبود عملکرد پرس می‌شود، در اغلب موارد به مستثنی کردن کل صفحات از جستجو، به ویژه برای پرس و جوهای دامنه منفی، اجازه می‌دهد.  
درخت  $R^*$  تلاش می‌کند برای کاهش هر دو؛ با استفاده از ترکیبی از یک الگوریتم تقسیم گره تجدید نظر و مفهوم درج مجدد اجباری از سرریز گره جلوگیری کند. این براساس مشاهدات ساختار درخت  $R$  که بسیار حساس به ترتیبی که در ورودی‌ها درج شده است، می‌باشد. بنابراین ساختار ساخت درج به احتمال زیاد زیر بهینه است. حذف و درج مجدد ورودی‌ها به پیدا کردن مکانی در درخت‌ها که مناسبتر از مکان اصلی خودشان می‌باشد، اجازه می‌دهد [۵]. در جدول ۲، مقایسه بین درخت  $R$  و  $R^*$  و  $R^+$  نشان داده شده است.

جدول ۲، مقایسه بین درخت  $R$  و درخت‌های  $R^*$  و  $R^+$

مقایسه	$R$
$R^+$	در $R^+$ گره‌های کمتری نسبت به $R$ مشاهده می‌شود.
$R^*$	ساختار ساخت درج به احتمال زیاد زیر بهینه است.

### ۶. نتایج

آزمایشات انجام شده نشان می‌دهد که  $STHist$  نسبت به متدهای پیشنهادی دیگر برتر است. با این وجود  $STHist$  پیچیدگی زمانی  $O(n^2)$  برای داده‌ی ۲ بعدی و  $O(n^3)$  برای داده‌ی ۳ بعدی را دارد. روش دیگر برای دستیابی به هیستوگرام فضایی تولید آن است که از ساختار شاخص گذاری فضایی مثل  $R_{tree}$  استفاده می‌کند.  
در این مقاله  $R^+_{tree}$  با  $R^*_{tree}$  و  $R_{tree}$  بعنوان شاخص متعارف از خانواده  $R_{tree}$  ها مورد مقایسه قرار داده شد.  $R^+_{tree}$  بعنوان نسخه اصلاحی  $R^+_{tree}$  ارائه گردید. نتیجه حاصل نشان داد که  $R^+_{tree}$  در دامنه‌ی پرس و جو، پرس و جو  $KNN$  و پرس و جو  $Top-k$  بسیار کارآمدتر از  $R^*_{tree}$  می‌باشد. کارآمدی  $R^+_{tree}$  با رشد ابعاد به آرامی کاهش می‌یابد، چون مقادیر تکرار هم افزایش می‌یابد. در انتها نشان داده شد که راندمان زمان جستجوی  $R^+_{tree}$  بسیار بهتر از  $R^+_{tree}$  است، آن هم زمانی که دیتاپوینت تکراری مد نظر قرار می‌گیرد می‌باشد.



## ۷. منابع

- [۱] Combining R-Tree and B-Tree to Enhance Spatial Queries Processing Marwa A. M. Abd Elwahab<sup>۱</sup>, Khaled M. Mahar<sup>۲</sup>, HatemAbdelkader<sup>۳</sup>, HatemAwad Khater<sup>۴</sup>
- [۲] New Database Architecture for Smart Query Handler of Spatial Database ParthasarathiBoyal, RituparnaChaki b aWest Bengal University of Technology, BF<sup>۴۲</sup> SaltLake City, Kolkata-۷۰۰۰۶۴, India bWest Bengal University of Technology, BF<sup>۴۲</sup> SaltLake City, Kolkata-۷۰۰۰۶۴, India
- [۳] A class of R-tree histograms for spatial databases Technical Report DaniarAchakeev Department of Mathematics and Computer Science Philipps-Universität Marburg, Germany [achakeye@mathematik.uni-marburg.de](mailto:achakeye@mathematik.uni-marburg.de) Bernhard Seeger Department of Mathematics and Computer Science Philipps-Universität Marburg, Germany [seeger@mathematik.uni-marburg.de](mailto:seeger@mathematik.uni-marburg.de)
- [۴] International Conference on Computational Science, ICCS ۲۰۱۲ Workshop on using Emerging Parallel Architectures (WEPA ۲۰۱۲) Speeding up spatial database query execution using GPUs Bogdan Simion<sup>۱,\*</sup>, Suprio Ray\*, Angela Demke Brown\*Department of Computer Science, University of Toronto
- [۵] R++-tree: an efficient spatial access method for highly redundant point data Martin Šumák, Peter Gurský P. J. Šafárik University in Košice, Jesenná ۵, ۰۴۰۰۱ Košice, Slovakia martin.sumak@student.upjs.sk, peter.gursky@upjs.sk
- [۶] Improving the Accuracy of Histograms for Geographic Data Objects ,HaiThanh Mai, Jaeho Kim, and Myoung Ho Kim  
Department of Computer Science, KAIST  
۲۹۱ Daehak - ro, Yu seong- Gu, D aeje on ۳۰۵- ۷۰۱, R epublic of Korea  
{mhthanh,jaeho,mhkim}@dbserver.kaist.ac.kr [۷] Y. Lifang; L. Rui; H. Xianglin; L. Yueping, "Performance of Rtree with slim-down & Reinsertion Algorithm," in ProcInternational Conference on Signal Acquisition and Processing, pages: ۲۹۱-۲۹۴, ۲۰۱۰.
- [۸] Z. Shaohui, C. Zhanwei , "The research of Hilbert R-tree spatial index algorithm based on Hybrid clustering," in Proc International Conference on Electronic & Mechanical Engineering and Information Technology, pages: ۳۴۹۵-۳۴۹۷, ۲۰۱۱.
- [۹] Kao, B. Lee, S. Lee, F. Cheung, D. " Clustering Uncertain Data Using Voronoi Diagrams and R-Tree Index,"IEEE Transactions on Knowledge and data engineering, Vol. ۲۲No. ۹, ۲۰۱۰.
- [۱۰] Shengnan, K. " Integrating R-tree and Levels of Detail,"Eighth International Conference on Fuzzy Systems and KnowledgeDiscovery (FSKD), ۲۰۱۱.
- [۱۱] Yu,B. Kim,H. Choi,W. Kwon,D. , " Parallel Range Query Processing on R-Tree with Graphics Processing Unit," Ninth IEEE International Conference on Dependable, Autonomic and Secure Computing, ۲۰۱۱