

Hybrid Vector Library—From Memory Bound to Compute Bound with NVVM

Régis PORTALEZ — ALTIMESH — regis.portalez@altimesh.com — Florent DUGUET — ALTIMESH — florent.duguet@altimesh.com

MOTIVATION

Existing source code usually interleaves data management, error-checking, text processing and actual compute. On general purpose processors, this mixture of code tasks is not necessarily an issue, and performance levels are often satisfactory as is.

However, when trying to use GPU, this hybrid computing turns into a coding challenge. Each individual computing tasks does not show sufficient workload, and porting the whole application requires a significant investment in the software asset.

We propose an alternate approach with runtime compilation based on function calls on a compute library. Hybrid Vector Library operates on vectors, in a manner similar to BLAS level 1 routines, with other functions such as square root or exponential, or MKL routines. In essence, all operations are performed on a vector of values. We illustrate the performance results of this approach on a typical financial benchmark.

Existing solutions such as ArrayFire [5] do not allow custom device function to be called in the middle of a level 1 routines sequence. We address that issue by processing these functions at compile time.

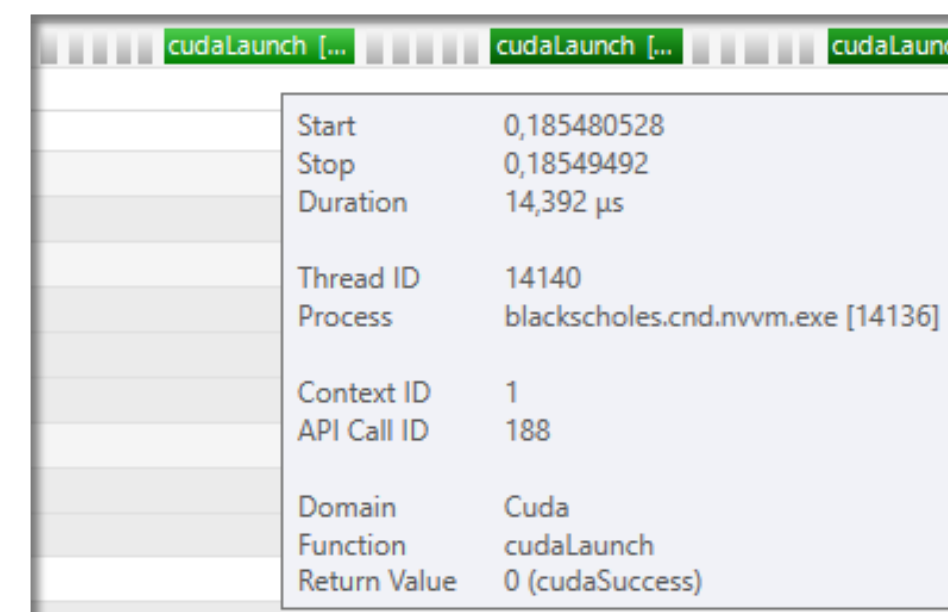
HYBRID VECTOR LIBRARY

Similar to MKL or BLAS Level-1 routines, Hybrid Vector Library exposes operations on vectors of values. These operations include basic arithmetic operations, along with mathematical function calls. It also exposes comparison tools and select operation to support basic value-dependent branching operations.

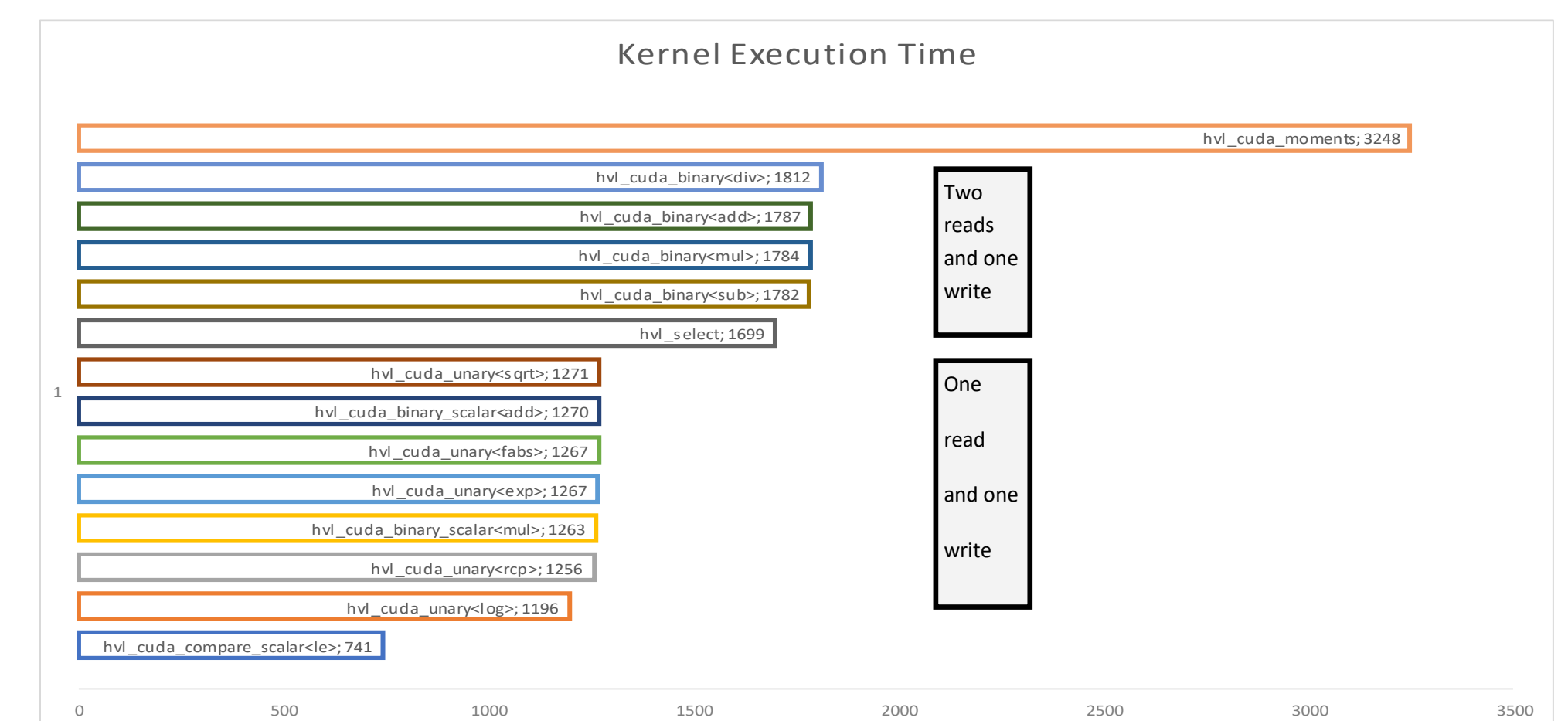
The API has several implementations that can be chosen at runtime to allow maximal flexibility. We illustrate here the use of two of these implementations.

PERFORMANCE OF NAIVE IMPLEMENTATION

The naive implementation will perform a kernel call for each vector operation. Beyond the lack of compiler optimization that would for example reconstruct FMA operations, this implementation suffers an important performance penalty. Indeed, each kernel call needs to be scheduled and executed. As illustrated in the following profiling snapshots, the execution time of a launch is about 25 microseconds (10µs configuration and 15µs launch). Within this time, about 1 million vector entries can be processed (calculating exp or log of the values for instance)



Moreover, kernel executions are memory bound. Indeed, current GPUs can execute more than 50 FLOPS for each memory operation, making all simple math functions, including transcendentals such as exponential, memory bound. We can see performance is driven by memory operations and not arithmetic complexity.



When executing operations upon library API call, performance is memory-bound and kernel execution time solely depends on amount of memory read or written.

RUNTIME COMPILATION

Depending on the implementation of HVL, execution of the calculation is performed at different stages. For the basic implementation, execution is done upon the API call on a vector of data. When using the NVVM-backed version, intermediate results do not exist. Operations are done in four phases:

1 User-defined device functions are identified in the call graph. CUDA source is generated for each of them, as long as a cubin file. The pairs function/cubin is registered at application startup.

2 When calling API methods, the operations are not scheduled immediately on the device. The different calls are gathered in a graph, which is by construction directed and acyclic (DAG), and no operation is executed until results are queried.

3 At given milestones, the DAG is converted into NVVM source code: each node is an NVVM statement with a single output. The NVVM source code is compiled at runtime. The sequence of calls and the compilation result are cached for future usage.

4 The resulting device binary is then scheduled for execution and results can be queried. The DAG is encoded into a signature in order to cache the compilation results — CUDA binary module. As shown, the compilation time may be longer than the overall execution time.

C++ Application Code (1)

```
// scalar code
double BlackScholesBodyScalar(
    double S, //Stock price
    double X, //Option strike
    double T, //Option years
    double R, //Riskless rate
    double V, //Volatility rate
)
{
    double S = S, X = X, T = T, R = R, V = V;

    double sqrtT = sqrt(T);
    double d1 = (log(S / X) + (R + 0.5 * V * V) * T) / (V * sqrtT);
    double d2 = d1 - V * sqrtT;
    double CND01 = CND(d1);
    double CND02 = CND(d2);

    double expRT = exp(-R * T);
    return (S * CND01 - X * expRT * CND02);
}

// vector code
// Earlier: mycnd has been declared extern "C" for symbol to be retrieved
hvlvect BlackScholesBodyVect(const hvlvect& S, const hvlvect& X, const hvlvect& T,
    const hvlvect& R, const hvlvect& V)
{
    hvlvect VsqrtT = V * sqrt(T);
    hvlvect d1 = (log(S / X) + (R + 0.5 * V * V) * T) / (VsqrtT);
    hvlvect d2 = d1 - VsqrtT;
    hvl_invoke(d1, mycnd);
    hvl_invoke(d2, mycnd); // (2)
    hvlvect expRT = exp(-R * T);
    return (S * d1 - X * expRT * d2);
}
```

(1) In this implementation all calls are within the same function but it can be spread on multiple source files or binary modules.



Vector operators are overloaded in C++ to make use of hvl library and include error-checking using exceptions. User defined device code is invoked through special API calls.

HVL API Calls

```
hvl_create
...
hvl_assign_hvlvect
hvl_apply_sqrt
hvl_assign_hvlvect
hvl_mul
hvl_assign_hvlvect
hvl_mul
hvl_mul_scalar
hvl_add
hvl_mul
hvl_assign_hvlvect
hvl_div
hvl_apply_log
hvl_add
hvl_div
hvl_assign_hvlvect
hvl_sub
...
hvl_invoke
hvl_invoke
hvl_assign_hvlvect
hvl_mul
hvl_mul_scalar
hvl_apply_exp
hvl_assign_hvlvect
hvl_mul
hvl_mul
hvl_sub
hvl_destroy
...
```

From HVL API Calls, a DAG is generated and upon request transformed into NVVM source code.

NVVM Generated Code

```
define void @hvl_nvmm_8 (i64 %n, double* %outptr, double* %param_load_37, double* %param_load_38, double* %param_load_41, double* %param_load_52, double* %param_load_54) {
@entry:
for.body:
  %load_idx_37 = getelementptr inbounds double* @param_load_37, i64 %idexptr
  %l_37 = load double* @load_idx_37, align 8
  %load_idx_38 = getelementptr inbounds double* @param_load_38, i64 %idexptr
  %l_38 = load double* @load_idx_38, align 8
  %l_39 = fmul double @l_37, %l_38
  %l_40 = call double @__my_log ( double %l_39 )
  %load_idx_41 = getelementptr inbounds double* @param_load_41, i64 %idexptr
  %l_41 = load double* @load_idx_41, align 8
  %load_idx_52 = getelementptr inbounds double* @param_load_52, i64 %idexptr
  %l_52 = load double* @load_idx_52, align 8
  %l_43 = fmul double %l_52, %l_41
  %l_44 = fmul double %l_43, %l_39
  %l_45 = fadd double %l_44, %l_40
  %load_idx_54 = getelementptr inbounds double* @param_load_54, i64 %idexptr
  %l_54 = load double* @load_idx_54, align 8
  %l_39 = fmul double %l_45, %l_54
  %l_55 = call double @__my_sqrt ( double %l_54 )
  %l_56 = fmul double %l_55, %l_55
  %l_33 = fdiv double %l_33, %l_56
  %l_4 = call double @mynd (double %l_33)
  %l_2 = fmul double %l_37, %l_4
  %l_28 = fmul double %l_41, -1.0000000000000000e+000
  %l_27 = fmul double %l_28, %l_54
  %l_26 = call double @__my_exp ( double %l_27 )
  %l_24 = fmul double %l_26, %l_26
  %l_32 = fsub double %l_31, %l_55
  %l_31 = call double @mynd (double %l_32)
  %l_23 = fmul double %l_24, %l_31
  %l_1 = fsub double %l_2, %l_23
  %outptr_idx = getelementptr inbounds double* %outptr, i64 %idexptr
  store double %l_1, double* %outptr_idx, align 8
  %idexptr.next = add i64 %idexptr, %stepprn
br label %for.body

for.body:
  ....
function_end:
```

A cubin file is generated at runtime and linked with precompiled modules of custom device functions.

SASS Generated Code

```
section ".nvvm.metadata"
align 4
hvl_nvmm_8
start_hvl_nvmm_8
end_hvl_nvmm_8
; entry block %v
; name compute block %v
; L1:
; L2:
; L3:
; L4:
; L5:
; L6:
; L7:
; L8:
; L9:
; L10:
; L11:
; L12:
; L13:
; L14:
; L15:
; L16:
; L17:
; L18:
; L19:
; L20:
; L21:
; L22:
; L23:
; L24:
; L25:
; L26:
; L27:
; L28:
; L29:
; L30:
; L31:
; L32:
; L33:
; L34:
; L35:
; L36:
; L37:
; L38:
; L39:
; L40:
; L41:
; L42:
; L43:
; L44:
; L45:
; L46:
; L47:
; L48:
; L49:
; L50:
; L51:
; L52:
; L53:
; L54:
; L55:
; L56:
; L57:
; L58:
; L59:
; L60:
; L61:
; L62:
; L63:
; L64:
; L65:
; L66:
; L67:
; L68:
; L69:
; L70:
; L71:
; L72:
; L73:
; L74:
; L75:
; L76:
; L77:
; L78:
; L79:
; L80:
; L81:
; L82:
; L83:
; L84:
; L85:
; L86:
; L87:
; L88:
; L89:
; L90:
; L91:
; L92:
; L93:
; L94:
; L95:
; L96:
; L97:
; L98:
; L99:
; L100:
; L101:
; L102:
; L103:
; L104:
; L105:
; L106:
; L107:
; L108:
; L109:
; L110:
; L111:
; L112:
; L113:
; L114:
; L115:
; L116:
; L117:
; L118:
; L119:
; L120:
; L121:
; L122:
; L123:
; L124:
; L125:
; L126:
; L127:
; L128:
; L129:
; L130:
; L131:
; L132:
; L133:
; L134:
; L135:
; L136:
; L137:
; L138:
; L139:
; L140:
; L141:
; L142:
; L143:
; L144:
; L145:
; L146:
; L147:
; L148:
; L149:
; L150:
; L151:
; L152:
; L153:
; L154:
; L155:
; L156:
; L157:
; L158:
; L159:
; L160:
; L161:
; L162:
; L163:
; L164:
; L165:
; L166:
; L167:
; L168:
; L169:
; L170:
; L171:
; L172:
; L173:
; L174:
; L175:
; L176:
; L177:
; L178:
; L179:
; L180:
; L181:
; L182:
; L183:
; L184:
; L185:
; L186:
; L187:
; L188:
; L189:
; L190:
; L191:
; L192:
; L193:
; L194:
; L195:
; L196:
; L197:
; L198:
; L199:
; L200:
; L201:
; L202:
; L203:
; L204:
; L205:
; L206:
; L207:
; L208:
; L209:
; L210:
; L211:
; L212:
; L213:
; L214:
; L215:
; L216:
; L217:
; L218:
; L219:
; L220:
; L221:
; L222:
; L223:
; L224:
; L225:
; L226:
; L227:
; L228:
; L229:
; L230:
; L231:
; L232:
; L233:
; L234:
; L235:
; L236:
; L237:
; L238:
; L239:
; L240:
; L241:
; L242:
; L243:
; L244:
; L245:
; L246:
; L247:
; L248:
; L249:
; L250:
; L251:
; L252:
; L253:
; L254:
; L255:
; L256:
; L257:
; L258:
; L259:
; L260:
; L261:
; L262:
; L263:
; L264:
; L265:
; L266:
; L267:
; L268:
; L269:
; L270:
; L271:
; L272:
; L273:
; L274:
; L275:
; L276:
; L277:
; L278:
; L279:
; L280:
; L281:
; L282:
; L283:
; L284:
; L285:
; L286:
; L287:
; L288:
; L289:
; L290:
; L291:
; L292:
; L293:
; L294:
; L295:
; L296:
; L297:
; L298:
; L299:
; L300:
; L301:
; L302:
; L303:
; L304:
; L305:
; L306:
; L307:
; L308:
; L309:
; L310:
; L311:
; L312:
; L313:
; L314:
; L315:
; L316:
; L317:
; L318:
; L319:
; L320:
; L321:
; L322:
; L323:
; L324:
; L325:
; L326:
; L327:
; L328:
; L329:
; L330:
; L331:
; L332:
; L333:
; L334:
; L335:
; L336:
; L337:
; L338:
; L339:
; L340:
; L341:
; L342:
; L343:
; L344:
; L345:
; L346:
; L347:
; L348:
; L349:
; L350:
; L351:
; L352:
; L353:
; L354:
; L355:
; L356:
; L357:
; L358:
; L359:
; L360:
; L361:
; L362:
; L363:
; L364:
; L365:
; L366:
; L367:
; L368:
; L369:
; L370:
; L371:
; L372:
; L373:
; L374:
; L375:
; L376:
; L377:
; L378:
; L379:
; L380:
; L381:
; L382:
; L383:
; L384:
; L385:
; L386:
; L387:
; L388:
; L389:
; L390:
; L391:
; L392:
; L393:
; L394:
; L395:
; L396:
; L397:
; L398:
; L399:
; L400:
; L401:
; L402:
; L403:
; L404:
; L405:
; L406:
; L407:
; L408:
; L409:
; L410:
; L411:
; L412:
; L413:
; L414:
; L415:
; L416:
; L417:
; L418:
; L419:
; L420:
; L421:
; L422:
; L423:
; L424:
; L425:
; L426:
; L427:
; L428:
; L429:
; L430:
; L431:
; L432:
; L433:
; L434:
; L435:
; L436:
; L437:
; L438:
; L439:
; L440:
; L441:
; L442:
; L443:
; L444:
; L445:
; L446:
; L447:
; L448:
; L449:
; L450:
; L451:
; L452:
; L453:
; L454:
; L455:
; L456:
; L457:
; L458:
; L459:
; L460:
; L461:
; L462:
; L463:
; L464:
; L465:
; L466:
; L467:
; L468:
; L469:
; L470:
; L471:
; L472:
; L473:
; L474:
; L475:
; L476:
; L477:
; L478:
; L479:
; L480:
; L481:
; L482:
; L483:
; L484:
; L485:
; L486:
; L487:
; L488:
; L489:
; L490:
; L491:
; L492:
; L493:
; L494:
; L495:
; L496:
; L497:
; L498:
; L499:
; L500:
; L501:
; L502:
; L503:
; L504:
; L505:
; L506:
; L507:
; L508:
; L509:
; L510:
; L511:
; L512:
; L513:
; L514:
; L515:
; L516:
; L517:
; L518:
; L519:
; L520:
; L521:
; L522:
; L523:
; L524:
; L525:
; L526:
; L527:
; L528:
; L529:
; L530:
; L531:
; L532:
; L533:
; L534:
; L535:
; L536:
; L537:
; L538:
; L539:
; L540:
; L541:
; L542:
; L543:
; L544:
; L545:
; L546:
; L547:
; L548:
; L549:
; L550:
; L551:
; L552:
; L553:
; L554:
; L555:
; L556:
; L557:
; L558:
; L559:
; L560:
; L561:
; L562:
; L563:
; L564:
; L565:
; L566:
; L567:
; L568:
; L569:
; L570:
; L571:
; L572:
; L573:
; L574:
; L575:
; L576:
; L577:
; L578:
; L579:
; L580:
; L581:
; L582:
; L583:
; L584:
; L585:
; L586:
; L587:
; L588:
; L589:
; L590:
; L591:
; L592:
; L593:
; L594:
; L595:
; L596:
; L597:
; L598:
; L599:
; L600:
; L601:
; L602:
; L603:
; L604:
; L605:
; L606:
; L607:
; L608:
; L609:
; L610:
; L611:
; L612:
; L613:
; L614:
; L615:
; L616:
; L617:
; L618:
; L619:
; L620:
; L621:
; L622:
; L623:
; L624:
; L625:
; L626:
; L627:
; L628:
; L629:
; L630:
; L631:
; L632:
; L633:
; L634:
; L635:
; L636:
; L637:
; L638:
; L639:
; L640:
; L641:
; L642:
; L643:
; L644:
; L645:
; L646:
; L647:
; L648:
; L649:
; L650:
; L651:
; L652:
; L653:
; L654:
; L655:
; L656:
; L657:
; L658:
; L659:
; L660:
; L661:
; L662:
; L663:
; L664:
; L665:
; L666:
; L667:
; L668:
; L669:
; L670:
; L671:
; L672:
; L673:
; L674:
; L675:
; L676:
; L677:
; L678:
; L679:
; L680:
; L681:
; L682:
; L683:
; L684:
; L685:
; L686:
; L687:
; L688:
; L689:
; L690:
; L691:
; L692:
; L693:
; L694:
; L695:
; L696:
; L697:
; L698:
; L699:
; L700:
; L701:
; L702:
; L703:
; L704:
; L705:
; L706:
; L707:
; L708:
; L709:
; L710:
; L711:
; L712:
; L713:
; L714:
; L715:
; L716:
; L717:
; L718:
; L719:
; L720:
; L721:
; L722:
; L723:
; L724:
; L725:
; L726:
; L727:
; L728:
; L729:
; L730:
; L731:
; L732:
; L733:
; L734:
; L735:
; L736:
; L737:
; L738:
; L739:
; L740:
; L741:
; L742:
; L743:
; L744:
; L745:
; L746:
; L747:
; L748:
; L749:
; L750:
; L751:
; L752:
; L753:
; L754:
; L755:
; L756:
; L757:
; L758:
; L759:
; L760:
; L761:
; L762:
; L763:
; L764:
; L765:
; L766:
; L767:
; L768:
; L769:
; L770:
; L771:
; L772:
; L773:
; L774:
; L775:
; L776:
; L777:
; L778:
; L779:
; L780:
; L781:
; L782:
; L783:
; L784:
; L785:
; L786:
; L787:
; L788:
; L789:
; L790:
; L791:
; L792:
; L793:
; L794:
; L795:
; L796:
; L797:
; L798:
; L799:
; L800:
; L801:
; L802:
; L803:
; L804:
; L805:
; L806:
; L807:
; L808:
; L809:
; L810:
; L811:
; L812:
; L813:
; L814:
; L815:
; L816:
; L817:
; L818:
; L819:
; L820:
; L821:
; L822:
; L823:
; L824:
; L825:
; L826:
; L827:
; L828:
; L829:
; L830:
; L831:
; L832:
; L833:
; L834:
; L835:
; L836:
; L837:
; L838:
; L839:
; L840:
; L841:
; L842:
; L843:
; L844:
; L845:
; L846:
; L847:
; L848:
; L849:
; L850:
; L851:
; L852:
; L853:
; L854:
; L855:
; L856:
; L857:
; L858:
; L859:
; L860:
; L861:
; L862:
; L863:
; L864:
; L865:
; L866:
; L867:
; L868:
; L869:
; L870:
; L871:
; L872:
; L873:
; L874:
; L875:
; L876:
; L877:
; L878:
; L879:
; L880:
; L881:
; L882:
; L883:
; L884:
; L885:
; L886:
; L887:
; L888:
; L889:
; L890:
; L891:
; L892:
; L893:
; L894:
; L895:
; L896:
; L897:
; L898:
; L899:
; L900:
; L901:
; L902:
; L903:
; L904:
; L905:
; L906:
; L907:
; L908:
; L909:
; L910:
; L911:
; L912:
; L913:
; L914:
; L915:
; L916:
; L917:
; L918:
; L919:
; L920:
; L921:
; L922:
; L923:
; L924:
; L925:
; L926:
; L927:
; L928:
; L929:
; L930:
; L931:
; L932:
; L933:
; L934:
; L935:
; L936:
; L937:
; L938:
; L939:
; L940:
; L941:
; L942:
; L943:
; L944:
; L945:
; L946:
; L947:
; L948:
; L949:
; L950:
; L951:
; L952:
; L953:
; L954:
; L955:
; L956:
; L957:
; L958:
; L959:
; L960:
; L961:
; L962:
; L963:
; L964:
; L965:
; L966:
; L967:
; L968:
; L969:
; L970:
; L971:
; L972:
; L973:
; L974:
; L975:
; L976:
; L977:
; L978:
; L979:
; L980:
; L981:
; L982:
; L983:
; L984:
; L985:
; L986:
; L987:
; L988:
; L989:
; L990:
; L991:
; L992:
; L993:
; L994:
; L995:
; L996:
; L997:
; L998:
; L999:
; L1000:
; L1001:
; L1002:
; L1003:
; L1004:
; L1005:
; L1006:
; L1007:
; L1008:
; L1009:
; L1010:
; L1011:
; L1012:
; L1013:
; L1014:
; L1015:
; L1016:
; L1017:
; L1018:
; L1019:
; L1020:
; L1021:
; L1022:
; L1023:
; L1024:
; L1025:
; L1026:
; L1027:
; L1028:
; L1029:
; L1030:
; L1031:
; L1032:
; L1033:
; L1034:
; L1035:
; L1036:
; L1037:
; L1038:
; L1039:
; L1040:
; L1041:
; L1042:
; L1043:
; L1044:
; L1045:
; L1046:
; L1047:
; L1048:
; L1049:
; L1050:
; L1051:
; L1052:
; L1053:
; L1054:
; L1055:
; L1056:
; L1057:
; L1058:
; L1059:
; L1060:
; L1061:
; L1062:
; L1063:
; L1064:
; L1065:
; L1066:
; L1067:
; L1068:
; L1069:
; L1070:
; L1071:
; L1072:
; L1073:
; L1074:
; L1075:
; L1076:
; L1077:
; L1078:
; L1079:
; L1080:
; L1081:
; L1082:
; L1083:
; L1084:
; L1085:
; L1086:
; L1087:
; L1088:
; L1089:
; L1090:
; L1091:
; L1092:
; L1093:
; L1094:
; L1095:
; L1096:
; L1097:
; L1098:
; L1099:
; L1100:
; L1101:
; L1102:
; L1103:
; L1104:
; L1105:
; L1106:
; L1107:
; L1108:
; L1109:
; L1110:
; L1111:
; L1112:
; L1113:
; L1114:
; L1115:
; L1116:
; L1117:
; L1118:
; L1119:
; L1120:
; L1121:
; L1122:
; L1123:
; L1124:
; L1125:
; L1126:
; L1127:
; L1128:
; L1129:
; L1130:
; L1131:
; L1132:
; L1133:
; L1134:
; L1135:
; L1136:
; L1137:
; L1138:
; L1139:
; L1140:
; L1141:
; L1142:
; L1143:
; L1144:
; L1145:
; L1146:
; L1147:
; L1148:
; L1149:
; L1150:
; L1151:
; L1152:
; L1153:
; L1154:
; L1155:
; L1156:
; L1157:
; L1158:
; L1159:
; L1160:
; L1161:
; L1162:
; L1163:
; L1164:
; L1165:
; L1166:
; L1167:
; L1168:
; L1169:
; L1170:
; L1171:
; L1172:
; L1173:
; L1174:
; L1175:
; L1176:
; L1177:
; L1178:
; L1179:
; L1180:
; L1181:
; L1182:
; L1183:
; L1184:
; L1185:
; L1186:
; L1187:
; L1188:
; L1189:
; L1190:
; L1191:
; L1192:
; L1193:
; L1194:
; L1195:
; L1196:
; L1197:
; L1198:
; L1199:
; L1200:
; L1201:
; L1202:
; L1203:
; L1204:
; L1205:
; L1206:
; L1207:
; L1208:
; L1209:
; L1210:
; L1211:
; L1212:
; L1213:
; L1214:
; L1215:
; L1216:
; L1217:
; L1
```