# 4.8 — Floating point numbers

👤 ALEX   📅 AUGUST 11, 2021

| Category | Type | Minimum Size | Typical Size |
|---|---|---|---|
| floating point | float | 4 bytes | 4 bytes |
| | double | 8 bytes | 8 bytes |
| | long double | 8 bytes | 8, 12, or 16 bytes |

```cpp
float fValue;
double dValue;
long double ldValue;
```

```cpp
int x{5}; // 5 means integer
double y{5.0}; // 5.0 is a floating point literal (no suffix means double type by
default)
float z{5.0f}; // 5.0 is a floating point literal, f suffix means float type
```

Mmrcęf_rę`wbcd_sjrñejm_rd eęnmd rę́grcp_jqdbcd_sjręmęrwncębms `jcñ@l ęąs dbyęrę qcbęrnbcl mrcę_ęgrcp_jęmdrwncędm_rñ

## Printing floating point numbers

Mmu ęaml qdbcpęf gęęck njcęnmrep_k ìJ

```cpp
#include <iostream>

int main()
{
  std::cout << 5.0 << '\n';
  std::cout << 6.7f << '\n';
  std::cout << 9876543.21 <<
'\n';

  return 0;
}
```

Sf cęcqs jrqęmdęf gęęck dj ejwęęk njcęnmrep_k ̀ęk_ węs pnpgcęwms ìJ

```
5
6.7
9.87654e+06
```

Herfcedpqre_qcĥeyfcerbIĴansrepojrcbeĵ ŀçtcl eʃmsefeuceʃvncbeĵ ɤ ffitAwebcd_sjrĥ qrbIĴansreugjeĵ mrenpoj reyfceb_argnl_jen_premte, I s k `cpeĝeyfceb_argnl_jen_preĝnĩ Ħ

Herfceꞯcamlbeꞯ_qcĥeyfceĵ s k `cpenpojrqe̜qeuceꞯvncarĦ

Herfcefgbbe_qcĥegrenpojrcbeyfceĵ s k `cpeĵ eꞯagcl rggbaĵ mr_rgnl ĝĵbvnseĵ ccbeꞯqcdpcqʃcpenl eꞯagcl rggbaĵ mr_rgnl Ḥqccėꞯcqqmlĭtĥ Herpmbsargnl eꞯmeꞯagcl rggbaĵ mr_rgnl ĝĦ

## Floating point range

@qqsk ĝ eeḦDDDeĵ ɩ Ĭ epcnpcqcl r_rgnl IJ

| Size | Range | Precision |
|---|---|---|
| Ĭ ęˋwrcq | ®Ĭ ffĬ e✓eĵ Ĩ ᶠ Ĭ ̇ eꞯme®Ĭ ffĬ e✓eĵ Ĩ ᵀ Ĭ ̇ | Ĭ ̂ eꞯe l ggba_l rebegrqĥrwnga_jjwĵ |
| Ĭ ęˋwrcq | ®Ĭ ffĬ e✓eĵ Ĩ ᶠ Ĭ ̇ eꞯme®Ĭ ffĬ e✓eĵ Ĩ ᵀ Ĭ ̇ | ĩ Ĭ ̂ Ĭ eꞯe l ggba_l rebegrqĥrwnga_jjwĵ Ĭ |
| ĩ Ĭ ̂ grqⱺrwnga_jjwsꞯqcqĵ Ĩ empĵ Ĭ ęˋwrcq͛ | ®Ĭ ffĬ e✓eĵ Ĩ ᶠ ⸱ ᵀ Ĭ ̇ eꞯme®Ĭ ffĬ e✓eĵ Ĩ ᵀ ⸱ ᵀ Ĭ ̇ | ĩ Ĭ ̂ Ĭ eꞯe l ggba_l rebegrq |
| ĩ Ĭ ęˋwrcq | ®Ĭ ffĬ e✓eĵ Ĩ ᶠ ⸱ ᵀ Ĭ ̇ eꞯme®Ĭ ffĬ e✓eĵ Ĩ ᵀ ⸱ ᵀ Ĭ ̇ | ĩ Ĭ ̂ Ĭ eꞯe l ggba_l rebegrq |

Sfcẽĩh̃ gẹdm_rḍeẹnmḍrẹywncẹgẹ,ẽ gẹmẹ,ẽfgẹrmpẹ,jẹ lmk _jwFNl ẹk mbcpḷenpmacẹqmpẹfẹgẹqẹrwngẹ_jjwẹgk njck cl rcbẹsẹqḍeẽĩempẽĩẹ̀wrcq ẽufgẹfẹgẹ,ẽk mpcẹ_rsp_jẹqẹcẹhpẹnpmacẹqmpẹrmf_lbjcg̃H

H̃ẹk _wẹcck ẹ,ẽgrjcẹmbbẹf_rẹfcẽĩh̃ gẹdm_rḍeẹnmḍrẹywncẹf_qẹfcẹ_kcẹ_lecẹ,ẹqẹfcẽĩh̃ wrcẹdm_rḍeẹnmḍrẹywncF̃Sfgẹẹgẹ̀ca_sqcẹfcw f_tcẹfcẹ_kcẹsk `cpẹmḍẹgẹcbgẹ_rcbẹmẹfcẹcvnmlcl rẽf̃ẹmuctcpĥẹfcẽĩh̃ wrcẹ sk `cpẹ_lẹqrmpcẹk mpcẹẹgẹl gbẹ_lrẹgẹgrqH̃

## Floating point precision

Bml qbcpẹfcẹẹb_argml ẽĩh̃ tF̃Sfcẹbcạk _jẹqcnpcqcl r_rgml emḍẹfgẹl sk `cpẹẹẽĩīīīīīīīīīīīīīīñ ẹugfẽĩNẹẹmḍeẹmsrẹmḍ dḍgrwF̃H̃ẹwnsẹẹcpc uprḍeẹfgẹl sk `cpml ẹẹngcacṃdn_ncpĥẹwnspẹ,ẹpk ẹumsjbẹcẹrẹgẹcbẹ_rẹmk cẹmḍ rĥẹlbẹwnsN̂bẹctclrs_jjwẹrmẹuprḍeṽ@lbẹfcẹl sk `cp wnsẹucpcẹ,ẹcḍrẹeẹumsjbẹ`cẹjmqcẹrmẽĩīīīīīīīīīñ ĥẹugfẽĩNẹẹmḍeẹmsrẹmḍ dḍgrwg̃̀srḍ mrcv_arjwH̃

Nl ẹẹmk nsrcpĥẹlẹgḍ grcẹcl erfẹlsk `cpẹumsjbẹcẹsogẹ,cḍ dḍgrcẹk ck mpwẹrmpcĥẹlbẹrwng_jjwẹucml jwẹf_tcẽĩempẽĩẹ̀wrcqF̃Sfgẹẹgk grcb k ck mpwẹk c_lqẹdm_rḍeẹnmḍrẹlsk `cpẹẹ_lẹml jwẹrmpcẹ,ẹacpr_gḍ ẹlsk `cpẹqẹ,ẹẹl gbẹ_lrẹgẹgrqĥẹlbẹf_rẹlwẹ,bbgml _jẹẹl gbẹ_lr bgẹgrqẹ,pcẹmḍrF̃Sfcẹl sk `cpẹf_rẹgḍ ẹars_jjwẹrmpcbẹujjẹcẹjmqcẹmẹfcẹcqqcbẹl sk `cpĥẹ̀srḍ mrcv_rH̃

Sfcẹ precision emḍeẹdm_rḍeẹnmḍrẹl sk `cpẹbcḍl qcḹfmuẹk _lwẹ dgḍg/g̃b_lrẹdgḍgrqẹgre_lẹcpcnpcqcl rẹugfṃsrḍ dmpk _rgml émqqH̃

V f clẹmsrns rrḍẹeḍẹdm_rḍeẹnmḍrẹl sk `cpĥẹqĥ̀rbĬjẹmsrẹ_qẹẹbcd_sjrẹnpcaggml emḍt.ĥẹfẹ_rẹẹgĥ̀ qq sk cẹjjẹdm_rḍeẹnmḍrẹt_ạg `jcqẹ,pc ml jwẹẹl gbẹ_lrẹmḍ bgẹgrqẹfcḱ dḍ qk ẹẹnpcaggml emḍrĥẹlbẹfl acẹgẹjjẹrpsl a_rcẹl wfḍeẹe,ḍrcpẹf_rH̃

Sfcẹhjjmwḍeẹnmrep_k ẹẹfmuẹqĥ̀rbĬjẹmsrẹrps a_rḍeẹemḍ bgẹgrqIJ

```
1   #include <iostream>

2   int main()
3   {
4       std::cout << 9.87654321f << '\n';
5       std::cout << 987.654321f << '\n';
        std::cout << 987654.321f << '\n';
        std::cout << 9876543.21f << '\n';
6       std::cout << 0.0000987654321f <<
    '\n';

7       return 0;
    }
```

Sf gẹ,ẹnmrep_k ẹ,srns rsqIJ

```
9.87654
987.654
987654
9.87654e+006
9.87654e-005
```

Mmrcẹf_rẹ_afẹmḍfcqẹmḍ jwẹf_tcẹẹl gbẹ_lrẹbgẹgrqH̃

@jqmẹmpcẹf_rẹqĥ̀rbĬjẹmsrẹẹujjẹu gr_fẹmḍns rrḍ eẹẹ sk `cpḍ ẹ,ẹagḍ rgbẹ,mr rgmḍ ẹẹmk cẹ,qcqF̃Ccncl bḍ eẹmḹ ẹfcẹmk nẹgcpĥẹfc cvnmlcl rẹujjẹwngẹ̀cẹ,bbcbẹ,ẹẹ gḍ sk `cpẹẹbgẹgrqF̃Ec_pḹ mḍ ẽĩĵĩĬcg̃ĩĬḹeẹẹfcẹ_kcẹ,ẹqĥ ẹĵĩĬcĬN̂sqẹ,rsẹgẹk c n_bbgḍ eẽĩMF̃Sfcẹ,ḍ sk cĥ̀rẹẹl sk `cpẹbgẹgrqẹbmḍj_wḍbẹ,ẹmk nḍpcagggḍ aĥ,Ugs_jẹꞏrsbmgẹsqẽĩqk cẹmḹfcpẹqẽ,ḍcẹ,qcpcẹfc Bꞏ ẹqrḍ lb_pbg̃H

Sfce sk `cpemdpgegrqemdnpcaggml e ejm rg eemmg ret_pg `jcef_qbcncl bqeml e`mfefcegxce jm rqef_tce jcqqe npcaggml ef_l ebms `jcqje lb rfcep_prgasj_pet_jsce`cg eeprmpcbeGmk cet_jscqef_tce mpcenpcaggml ef_l emfcpqFEjm ret_jscqef_tce`cru ccl ere l be ebgegqemd npcaggml êu grf k mqrejm ret_jscqef_tg ee rejc_qrej eqe l gba_l rebgegqFCms `jcet_jscqef_tce`cru ccl êï e l be Ï ebgegqemdnpcaggml êu grf k mqrebms `jcet_jscqef_tg ee rejc_qrej eqe l gba_l rebgegqFKml ebms `jcef_qe ek g jk sk enpcaggml emdï Ñçï Ñmpe ï eqe l gba_l rebgegrq bcncl bg eeml ef mu ek _l we`wrcqegrmaas ngcqF
 
V cqe_l emtcppgbcefcebcd sjrenpcaggml ef_rqrblIjms rqef mu qe`wsqg ee l output manipulator eblargml e l _k cb std::setprecision() H
Output manipulatorsejrcpef mu  b_r_egems rnsrĥe l be pcebcdj cbej efcgnk_ l gnefc_bcpH

```cpp
1   #include <iostream>
2   #include <iomanip> // for output manipulator std::setprecision()

3   int main()
4   {
5       std::cout << std::setprecision(16); // show 16 digits of precision
6       std::cout << 3.33333333333333333333333333333333333333f <<'\n'; // f suffix means float
7       std::cout << 3.33333333333333333333333333333333333333333 << '\n'; // no suffix means double

8       return 0;
    }
```

Ns rns rqIJ

```
3.333333253860474
3.333333333333334
```

Aca_sqceu cqcrefcenpcaggml emêï ebgegqes qg e std::setprecision() ĥc_afmdrfce `mtce sk `cpqgegmpg rcbeu gfêï ebgegrqFAsrĥe q wns e_l eqccĥrfce l sk `cpqqecpr_g jwe pcl Nenpcaggcrmeê ï ebgegqê@l be ca_sqc djm rqe pcejcqqenpcagcef_l ebms `jcqĥrfcedjm rf_qek mpc cppmpH

Opcaggml egxqscqbml Ns rqrek n_areq p_argml _jd sk `cpqĥrfcwek n_are l we sk `cpe u gfemmk _l we gl gba_l rebgegrqFKcrNqeml qdbcpe ê ge l sk `cpIJ

```cpp
#include <iomanip> // for std::setprecision()
#include <iostream>

int main()
{
    float f { 123456789.0f }; // f has 10 significant digits
    std::cout << std::setprecision(9); // to show 9 digits
    in f
    std::cout << f << '\n';

    return 0;
}
```

Ns rns rIJ

```
123456792
```

ïÏïĭ ĭĬ ¡¡ ı Ïĝ ğ epc_rcpęrf_l ęĭ Ĭ ïĬ ïĬ ¡¡ ı ı Ĭ ı ĝ ç ę _ js cęĭ Ĭ ï Ĭ ï Ĭ ¡¡ ı Ì ı ĭ ñ ę f_ ę çĭ ĭ ę ęg l ğ ba_l r ğ ę grğ ñ `s rę ĭ m_rę _ js cęŗw ng a_j j w ęf_t cę ğ ę gr ą ę md ęn pc ag ğ ml Ġ ĭ l bęf cę qc sĭ jŗ emdĭ ĭ Ïĭ ĭ Ïĭ ĭ Ïĝ ¡¡ ı ı Ĭ ĝ emp cag gcę ml ĭ jw ęr mę ę ğ el ğ ba_l r ğ ę grğ Ĭ ñ V cę ĭ mр ę mk cę np ca gg ml ěñ V fc l ęn pc ag gml ę ę ğ ęi mŗ ę ` ca_s qc ę ę l s k `c p a_l Ñ ę ` cę ŗ mpc bę np ca gc jwĭ ęr f ğ ę ğ ę _ jj cb ę ę rounding error ñ

Bml qc os cl rj wñ ml cę f_ ğ ęr mę ` ca_ pc ds jµ f cl ę ğd ee ĵ m_r ğ ę ę m mğ rę l s k `c p ğ ęr f _ r ę pc os ğ ŗc ę mp cę np ca gg ml ęf_l ęr fc ę _ pğ `c qę _ l ę ĭ mjb ñ
```

> **Best practice**
>
> E_tmp ę ĥ ŕ ms `jc ęr ŕ cp ę ĵ m_r ęs l j cq ŗ ę ę n_a cę ŗ ę ęr p ck ĵs k ñ ę ęr fc ę _i emd ęn pc ag gml ęĝ ľ ęğ ĭ m_r ęu jj ęm ŕr cl ęĭ c_ bĭ mĵ l _ a as p_ agŗ cę ñ

## Rounding errors make floating point comparisons tricky

```
Ejm_rğ ę ę m mğ rę l s k `c pqc ęr pc ŕ pğ ŕ ęi wĵ mę u mpi ęu ĭ gr f ęĵ ĭ bs cę ŗm ĵ l ñm `t gr ms qę ĝ b gĭ Ĭ ĭ pc l ac ŗq ę `c ru cc l ę̀ ġ ĵ _ pŗw ĕ f mu ę ĭ b_ r_ ŕg ŗ ę mŗ pc bĭ ĝ ĵ l bĭ ęb ca ŗ ĝ k _ jⱥ ĭ f mu u c ę fğ ĭ ĵ ĝ ę l s k `c pŗ ĵ ŕ ĵ ĝ qb cŗ rf cę b_ arğ ml ęĭ ĭ fĭ ì Ĭ ĭ ñ Ĥ ĭ bc aĝ k _ jⱥ ĭ ęr f ğ ę mę _ qĵ ĝ w ęn pc qc qĭ l r cbĭ ę qĭ ŕ fìĭ ê ę l bĭ ęu cę pc s qc bĭ ęr mf ğ ĭ ĝ ŕ ę ę md ĭ ñ fŕ ę qĵ l c_ qĵ ĝ wĵ ĝ ęn pc qc qĭ l r_ `ĵ c l s k `c pⱥ ęu ĭ ğ ŗ ę ę l ğ ba_l r ğ ę grğ ñ Ġ mu c tc pŗ ñ ğ ę̀ ĝ ğ _ pŗw ĭ ñ ŕ fŕ ę ęr mę cę np cq c_l r cbĭ ę ŗw ęr fcĵ ĝ dĵ ĝ gr cę ęc os cl acĭ IJ ĭ ñ ĭ Ĭ ĭ ïï ĭ Ĭ ïï ĭ Ĭ ïï ĭ Ĭ ïï ĭ Ĭ ïï ĭ Ĭ ïï ĝ ŕ ñ ç Aca_s qc mŗ ęr f ĝ ĵ ĝ µ f cl ęu cę ŗ qĵ qŗ cl ę ĭ ñ f, ĭ ŗm ĝ ę ę m_r ğ ę ę m mğ rę l s k `c pⱥ ęu ĝ cl Ŋ ęĵ ĝ ęs l ęğ rm ĭ en pc ag gml ęn pm jck qĭ ñ
```

Xms ęĝ ĭ lęę qc cęr fc ęd b_ arq ęmd ęr f ğ ę µ ğ ęr f cę ĭ jj mu ğ ę e mp mep_k IJ

```cpp
#include <iomanip> // for std::setprecision()
#include <iostream>

int main()
{
    double d{0.1};
    std::cout << d << '\n'; // use default cout precision
    of 6
    std::cout << std::setprecision(17);
    std::cout << d << '\n';

    return 0;
}
```

Sf ğ ĝ ĝ ms rns rqĭ IJ

```
0.1
0.10000000000000001
```

NI ęf cęr mę ĵ ĝ c ĥ ŕ qr bĭ IJ ĝ ms rę pğ rqĭ ñ fŕ ę qĵ u cę cę vn ca r ñ

Nl efce`mrmk ej cÑu fcpce ucef_tcerblĴ msref mu e sqéį bgegroemden pcaggml Ñu cecce f_rbegę ars_jjwemros grcãflïeŠf gegę`ca_sqcefc bms_jcef_bemepsla_rcefce nnpmvgk _rgml bscemgrqęjck grcbęck mpwÑŠfcecosjrege elsk `cpef_regenpcaggcemęĮegel gbęa_lrbgegrq Ćufg afevwnceoms `jces_p Lrccqĝ'íĝ'srefcelsk `cpgml mrecv_arjwęĩflłĴQmsl bgęecppmpqęk_wqk_icęelsk `cpęgrfcpęjgefrjwęk_jjcpmp qjgefrjwęj_pecpÑ̂bcncl bgęeml eu fcpcefcepsla_rgml ef_nnclqĤ

Qmsl bgęecppmpqęa_lęf_tcęl cvncarcbemlqcosclacqĴĴ

```cpp
#include <iomanip> // for std::setprecision()
#include <iostream>

int main()
{
    std::cout << std::setprecision(17);

    double d1{ 1.0 };
    std::cout << d1 << '\n';

    double d2{ 0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1 }; // should equal 1.0
    std::cout << d2 << '\n';

    return 0;
}
```

```
1
0.99999999999999989
```

@jrfmsefęucęk gefręvncareęf_rbĩ ęl bębĨefmsjbę`cecos_jÑu cecce f_refcwępce l mrÑĤÑu ceu cpcemęamk n_pcbĩ ęl bębĨe ęl ęę npmep_k Ñefc npmep_k ęu msjbęnpm`_`jwęl mręcpdmpk ę_qęcvncarcbÑ̂Aca_sqcędjm_rg eęnmgl rę lsk `cpę ęrcl bemrę`cęcl cv_arÑemk n_pjl eędjm_rg eęnmgl r lsk `cpęgeencl cp_jjwęnpm`jck _rgæfÑu ce bqe ogasqqefcps `hcarek mpcęĮ lbęmjsrgml qĵgęjcqqml ęĩ̂fÑ̂Qcj_rgml _jemcp rmpqę l bedjm_rg eęnmgl r amk n_pgml qĤ

Nl cę_qręl mrcęml qmsl bgęecppmpqĴĴ_rfck rga_jemcp rgml qĢĝsafcęę_c bbgrml ęl bęk sjrgnjga_rgml ęęrcl bemrę_icęsl bgęecppmpqęepmuÑ̂ Rmęcrcl efmsefęĝfíĤ_ęsl bgęecppmpqęd efcęĵjrfęel gbęl rbgegrÑu fcl euę bbęĩ̂fęrcl emk cqÑrfcęsl bgęecppmpef_qępcnrędrmpfc ĩ̂Įrfęel gbęl rbgegrÑĤBml rdęsceb mp rgml qęu msjbę_ sqęfgęecppmpmpę`camk cdj apc_qdęejwęel gbęl rÑ̂

<div class="key-insight">

## Key insight

Qmsl bgęecppmpqęa aspęu fcl ę ęelsk `cpęl Ñę`cęrmpcbęn pcaggcjwÑ̂Šfgmę l ęf_nnclętcl ue fgeęk njcęels k `cpÑĝ'jgę'cãÑĤŠfcpcdmpęÑ pms l bgę ecppmpq_l ę l bębmÑ̂f_nnclęjjefcęk cÑ̂Qmsl bgęecppmpqępcl Ñefcevacrgml Ñ̂Ñ̂fcwÑępcefcsjcÑMctcpęqqsk cęumsp djm_rgeemgl rę lsk `cpęe_pcęv_arÑ̂

@ampmjj_pwemdfges jceglĴ̂`cęu_pwemdję eeljm_rg eęnmgl r lsk `cpędmpębęj_cgjempaspęlawę`_r_Ĥ

</div>

## NaN and Inf

Sfcpcę pcęu mencagje_rcempgrcomedjm_rg eęnmgl r lęsk `cpqĤSfcębqrmgępegeInfÑu fgaf cpcnpcq lrqędljg rwÑĤ Ŋda_lę`_cnmsgrqcempcę_rgĴcHŠfcęęcamlbępcęNaNÑu fgaf fęrlbęlbqędmpÑMmrre_eMs k `cpÑŌĴSfcpcępcęvcrę l̂dlcpl rejbqrmędqmdM_MÑ̂ufg fęu ceu ml Ñ̂Ñbgeas qecfcpcęÑ̂M_M _ nbĤ fcęlwwet_gj_cgrf cemk njcpsqcę egnca qhmędpmpihk _rgĝÑDDDęĩ̂Ĩgehmjm djm_rgeel selsk `cpÑ̂Ĥl mpfcpb mp_ren qqcbÑefc dljjmu g eęecmbcenrmbseacqsl bcdĵl bcę cf_tgpÑ̂Ĥ

Here's a program that shows these in action:

```
#include <iostream>

int main()
{
    double zero {0.0};
    double posinf { 5.0 / zero }; // positive infinity
    std::cout << posinf << '\n';

    double neginf { -5.0 / zero }; // negative infinity
    std::cout << neginf << '\n';

    double nan { zero / zero }; // not a number (mathematically
invalid)
    std::cout << nan << '\n';

    return 0;
}
```

And the results using Visual Studio 2022 on Windows:

```
1.#INF
-1.#INF
1.#IND
```

> **Best practice**
>
> @tmb¦bgtggml e`wÃ¡_jrmecrf cp¦tctcl eqwnspamk ngcp¡qsnnmprqegH

## Conclusion

Smqsk k_pgxcÌfceu mefg eqewnseqfmsjbeck ck `cp¡`msrqjm_rgeenmgr¡lsk `cpqIJ

Î¶qEjm_rg eenmgr¡lsk `cpq¡pces qcds jefmpgmpgeetcpwj_pecmpetcpweÄ_jj¡k `cpq̧g ajsbg eefmqceu grf¡b_argl_jamk nml clrqH
Î¶qEjm_rg eenmgr¡lsk `cpqmdrcl f_tcq¡k_jjmslbg eepmpmpÌctcl u fcl efc¡sk `cpf_qtu cpaqel gb_lrbegrq¡f_lefcpcacggmH
   L_wqk cq¡fcqceml lmgacb¦ca_sqcefcwe pc¡mpmk_jjÌlb ¦ca_sqcfcl sk `cpqpces a_rcbefmpmsrnsrḨGmu ctcpÂ
   amk n_pgmlqemdjm_rg eenmgr¡lsk `cpqk_wj mregefcevncarcbecqsjrqÄDcrfmpmg eek_rfck_ra_jemcp_rgmlqml¡fcqce¡_scq
   u gjje_sceefcqmslbg eepmpmpmepmu j_pecpH



### Mcvr jcqqml
🔼 Ammjc_l t_jscq



### A_ai rm _ `jcednslrclrq



### Opctgmsq jcqqml
🔼 H rpmbs argml eymeaglrgdä_mr_rgml

Leave a comment... Put C++ code between triple-backticks (markdown style):```Your C++ cod

Name*

@ Email*

@t_r_pqqbmk ę frrnqThep_t_r_phmk ħę pcęml l carcbęmęvns pęnpmt gocbęk _gę bbpcqqFʰ

Mmrglwłk cę `ms rępcnjg:qlę 🔔 ę     **POST COMMENT**

**444 COMMENTS**

Newest ▾

Ⓧ

Ⓧ