

مسئله در مورد تصحیح اشتباهات چاپی در متن بدون استفاده از دیکشنری است. به شما متن حاوی تعداد زیادی خطای چاپی داده می شود و هدف این است که تا حد امکان بسیاری از خطاها را اصلاح کنید.

در این مسئله وضعیت ها به حروف درستی که باید تایپ شود اشاره دارد و خروجی به حروف واقعی که تایپ شده اشاره دارد.

دنباله ای از خروجی ها داده شده (به عنوان مثال حروف تایپ شده واقعی) و مسئله ساخت دوباره دنباله حات های مخفی است

(دنباله ای از حروف مورد نظر). بنابراین داده ها برای مسئله زیر به شکل زیر خواهد بود:

i i
n n
t t
r r
o o
d x
u u
c c
t t
i i
o i
n n
- -
t t
h h
e e

که ستون چپ حروف درست و ستون سمت راست شامل متن با خطا است.

داده ها برای این مسدله به صورت زیر ایجاد شده. با یک فایل متنی شروع می کنیم. برای سادگی تمام اعداد و نشانه گذاری ها به فاصله سفید و تمام حروف به حالت حروف کوچک تبدیل شده است. متن باقی مانده دنباله ای از حروف کوچک و کارکتر فاصله است که با علامت - در فایل داده مشخص شده است. در ادامه علط های چاپی با صورت مصنوعی به صورت زیر به داده ها اضافه شده است:

با احتمال 90 درصد حرف نوشته می شود اما با احتمال 10 درصد، یک همسایه تصادفی (در در ترتیب فیزیکی صفحه کلید) از حرف به جای حرف درست نوشته می شود. کارکتر فاصله همیشه به صورت درست نوشته می شود. در یک نوع سخت تر از مسئله نرخ خطا به 20 درصد افزایش می یابد.

تقریباً 20000 کاراکتر اول فایل برای تست کنار گذاشته شده است. 161000 کارکتر بقیه برای آموزش استفاده می شود.

به عنوان مثال، فایل اصلی شروع می شود با :

introduction the industrial revolution and its consequences have been a
disaster for the human race they have greatly increased the life expectancy
of those of us who live in advanced countries but they have destabilized
society have made life unfulfilling have subjected human beings to
indignities have led to widespread psychological suffering in the third world
to physical suffering as well and have inflicted severe damage on the natural

world the continued development of technology will worsen the situation it will certainly subject human beings to greater indignities and inflict greater damage on the natural world it will probably lead to greater social disruption and psychological suffering and it may lead to increased physical suffering even in advanced countries the industrial technological system may survive or it may break down if it survives it may eventually achieve a low level of physical and psychological suffering but only after passing through a long and very painful period of adjustment and only at the cost of permanently reducing human beings and many other living organisms to engineered products and mere cogs in the social machine

با خطای 20 درصد، شبیه به متن پایین خواهد بود:

introduction the industrial revolution and its consequences have been a disaster for the human race they have greatly increased the life expectancy of those of us who live in advanced countries but they have festabupusee cocisty have made live interiorilling have wibjested human beings to incingitids have led to widespread psychological suffering in the third world to physical suffering as well and have inflicted severe damage on the natural world the continued development of technology will worsen the situation it will certainly subject human beings to greater indignities and inflict greater damage on the natural world it will probably lead to greater social disruption and psychological suffering and it may lead to increased physical suffering even in advanced countries the industrial technological system may survive or it may break down if it survives it may eventually achieve a low level of physical and psychological suffering but only after passing through a long and very painful period of adjustment and only at the cost of permanently reducing human beings and many other living organisms to engineered products and mere cogs in the social machine

نرخ خطا حدود 16.5 درصد است (کمتر از 20 درصد چون کاراکترهای فاصله درست تایپ شده است).

متن دوباره ساخته شده به وسیله HMM با الگوریتم ویتربی شبیه به متن زیر است:

introduction the industrial revolution and its consequences have been a disaster for the human race they have greatly increased the life expectancy of those of us who live in advanced countries but they have festabupusee cocisty have made live interiorilling have wibjested human beings to incingitids have led to widespread psychological suffering in the third world to physical suffering as well and have inflicted severe damage on the natural world the continued development of technology will worsen the situation it will certainly subject human beings to greater indignities and inflict greater damage on the natural world it will probably lead to greater social disruption and psychological suffering and it may lead to increased physical suffering even in advanced countries the industrial technological system may survive or it may break down if it survives it may eventually achieve a low level of physical and psychological suffering but only after passing through a long and very painful period of adjustment and only at the cost of permanently reducing human beings and many other living organisms to engineered products and mere cogs in the social machine

نرخ خطا به حدود 10.4 درصد کاهش یافت.

داده های مربوط به مقدار دهی در typos10.data و typos20.data که شامل داده ها با خطای به ترتیب 10 درصد و 20 درصد است.